



## EUROPEAN PATENT APPLICATION

(43) Date of publication:  
02.01.2004 Bulletin 2004/01

(51) Int Cl.7: G06F 17/60

(21) Application number: 03006814.2

(22) Date of filing: 26.03.2003

(84) Designated Contracting States:  
AT BE BG CH CY CZ DE DK EE ES FI FR GB GR  
HU IE IT LI LU MC NL PT RO SE SI SK TR  
Designated Extension States:  
AL LT LV MK

(72) Inventors:  
• Goodman, Joshua Theodore  
Redmond, Washington 98052 (US)  
• Rounthwaite, Robert L.  
Fall City, Washington 98024 (US)

(30) Priority: 26.06.2002 US 180565

(74) Representative: Grünecker, Kinkeldey,  
Stockmair & Schwanhäusser Anwaltssozietät  
Maximilianstrasse 58  
80538 München (DE)

(71) Applicant: MICROSOFT CORPORATION  
Redmond, Washington 98052-6399 (US)

## (54) SPAM detector with challenges

(57) A system and method facilitating detection of unsolicited e-mail message(s) with challenges is provided. The invention includes an e-mail component and a challenge component. The system can receive e-mail message(s) and associated probabilities that the e-mail message(s) are spam. Based, at least in part, upon the associated probability, the system can send a challenge

to a sender of an e-mail message. The challenge can be an embedded code, computational challenge, human challenge and/or micropayment request. Based, at least in part, upon a response to the challenge (or lack of response), the challenge component can modify the associated probability and/or delete the e-mail message.

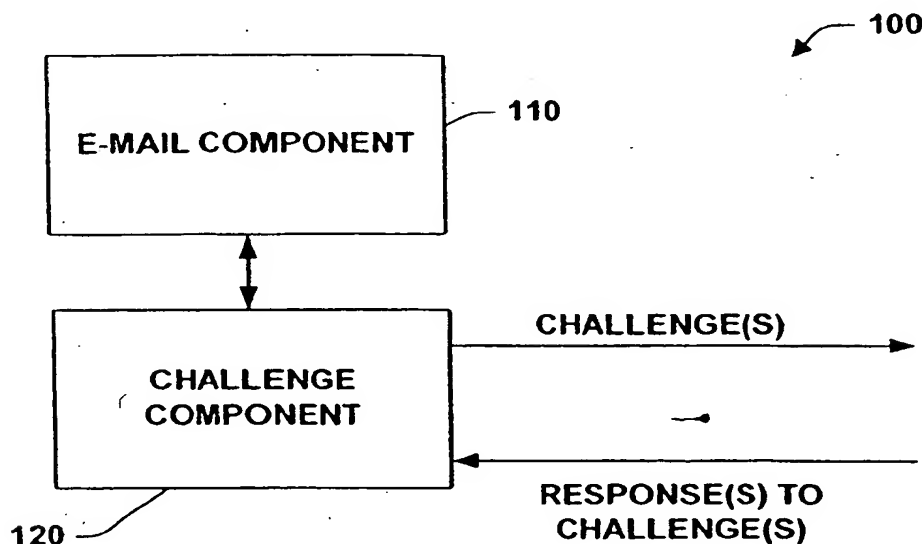


FIG. 1

## Description

### TECHNICAL FIELD

[0001] The present invention relates generally to electronic mail (e-mail) and more particularly to a system and method employing unsolicited e-mail (spam) detection with challenges.

### BACKGROUND OF THE INVENTION

[0002] Electronic messaging, particularly electronic mail ("e-mail") carried over the Internet, is rapidly becoming not only pervasive in society but also, given its informality, ease of use and low cost, a preferred mode of communication for many individuals and organizations.

[0003] Unfortunately, as has occurred with more traditional forms of communication (e.g., postal mail and telephone), e-mail recipients are increasingly being subjected to unsolicited mass mailings. With the explosion, particularly in the last few years, of Internet-based commerce, a wide and growing variety of electronic merchandisers is repeatedly sending unsolicited mail advertising their products and services to an ever expanding universe of e-mail recipients. Most consumers that order products or otherwise transact with a merchant over the Internet expect to and, in fact, regularly receive such merchant solicitations. However, electronic mailers are continually expanding their distribution lists to penetrate deeper into society in order to reach ever increasing numbers of recipients. For example, recipients who merely provide their e-mail addresses in response to perhaps innocuous appearing requests for visitor information generated by various web sites, often find, later upon receipt of unsolicited mail and much to their displeasure, that they have been included on electronic distribution lists. This occurs without the knowledge, let alone the assent, of the recipients. Moreover, as with postal direct mail lists, an electronic mailer will often disseminate its distribution list, whether by sale, lease or otherwise, to another such mailer, and so forth with subsequent mailers. Consequently, over time, e-mail recipients often find themselves barraged by unsolicited mail resulting from separate distribution lists maintained by a wide and increasing variety of mass mailers. Though certain avenues exist, based on mutual cooperation throughout the direct mail industry, through which an individual can request that his(her) name be removed from most direct mail postal lists, no such mechanism exists among electronic-mailers.

[0004] Once a recipient finds him(her)self on an electronic mailing list, that individual can not readily, if at all, remove his(her) address from it, thus effectively guaranteeing that he(her) will continue to receive unsolicited mail -- often in increasing amounts from that list and oftentimes other lists as well. This occurs simply because the sender either prevents a recipient of a message from

identifying the sender of that message (such as by sending mail through a proxy server) and hence precludes the recipient from contacting the sender in an attempt to be excluded from a distribution list, or simply ignores any request previously received from the recipient to be so excluded.

[0005] An individual can easily receive hundreds of unsolicited postal mail messages over the course of a year, or less. By contrast, given the ease and insignificant cost through which e-distribution lists can be readily exchanged and e-mail messages disseminated across large numbers of addressees, a single e-mail addressee included on several distribution lists can expect to receive a considerably larger number of unsolicited messages over a much shorter period of time. Furthermore, while many unsolicited e-mail messages (e.g., offers for discount office or computer supplies or invitations to attend conferences of one type or another) are benign; others, such as pornographic, inflammatory and abusive material, can be highly offensive to certain recipients.

[0006] Unsolicited e-mail messages are commonly referred to as "spam". Similar to the task of handling junk postal mail, an e-mail recipient must sift through his(her) incoming mail to remove spam. Unfortunately, the choice of whether a given e-mail message is spam or not is highly dependent on the particular recipient and content of the message - what may be spam to one recipient may not be so to another. Frequently, an electronic mailer will prepare a message such that its true content is not apparent from its subject line and can only be discerned from reading the body of the message. Hence, the recipient often has the unenviable task of reading through each and every message he(her) receives on any given day, rather than just scanning its subject line, to fully remove spam messages. Needless to say, such filtering (often manually-based) can be a laborious, time-consuming task.

[0007] In an effort to automate the task of detecting abusive newsgroup messages (so-called "flames"), the art teaches an approach of classifying newsgroup messages through a rule-based text classifier. See, E. Spertus "Smokey: Automatic Recognition of Hostile Messages", Proceedings of the Conference on Innovative Applications in Artificial Intelligence (IAAI), 1997. Here, semantic and syntactic textual classification features are first determined by feeding an appropriate corpus of newsgroup messages, as a training set, through a probabilistic decision tree generator. Given handcrafted classifications of each of these messages as being a "flame" or not, the generator delineates specific textual features that, if present or not in a message, can predict whether, as a rule, the message is a flame or not. Those features that correctly predict the nature of the message with a sufficiently high probability are then selected for subsequent use. Thereafter, to classify an incoming message, each sentence in that message is processed to yield a multi-element (e.g., 47 element) feature vector,

with each element simply signifying the presence or absence of a different feature in that sentence. The feature vectors of all sentences in the message are then summed to yield a message feature vector (for the entire message). The message feature vector is then evaluated through corresponding rules produced by the decision tree generator to assess, given a combination and number of features that are present or not in the entire message, whether that message is either a flame or not. For example, as one semantic feature, the author noticed that phrases having the word "you" modified by a certain noun phrase, such as "you people", "you bozos", "you flamers", tend to be insulting. An exception is the phrase "you guys" which, in use, is rarely insulting. Therefore, one feature is whether any of these former word phrases exist. The associated rule is that, if such a phrase exists, the sentence is insulting and the message is a flame. Another feature is the presence of the word "thank", "please" or phrasal constructs having the word "would" (as in: "Would you be willing to e-mail me your logo") but not the words "no thanks". If any such phrases or words are present (with the exception of "no thanks"), an associated rule, which the author refers to as the "politeness rule" categorizes the message as polite and hence not a flame. With some exceptions, the rules used in this approach are not site-specific, that is, for the most part they use the same features and operate in the same manner regardless of the addressee being mailed.

[0008] A rule based textual e-mail classifier, here specifically one involving learned "keyword-spotting rules", is described in W. W. Cohen, "Learning Rules that Classify E-mail", 1996 AAAI Spring Symposium on Machine Learning in Information Access, 1996 (hereinafter the "Cohen" publication). In this approach, a set of e-mail messages previously classified into different categories is provided as input to the system. Rules are then learned from this set in order to classify incoming e-mail messages into the various categories. While this method does involve a learning component that allows for automatic generation of rules, these rules simply make yes/no distinctions for classification of e-mail messages into different categories without providing any confidence measure for a given prediction. Moreover, in this work, the actual problem of spam detection was not addressed. In this regard, rule-based classifiers suffer various serious deficiencies which, in practice, would severely limit their use in spam detection. First, existing spam detection systems require users to manually construct appropriate rules to distinguish between legitimate mail and spam. Most recipients will not bother to undertake such laborious tasks. As noted above, an assessment of whether a particular e-mail message is spam or not can be rather subjective with its recipient. What is spam to one recipient may, for another, not be. Furthermore, non-spam mail varies significantly from person to person. Therefore, for a rule based-classifier to exhibit acceptable performance in filtering most spam

from an incoming mail stream, the recipient must construct and program a set of classification rules that accurately distinguishes between what constitutes spam and what constitutes non-spam (legitimate) e-mail. Properly doing so can be an extremely complex, tedious and time-consuming task even for a highly experienced and knowledgeable computer user.

[0009] Second, the characteristics of spam and non-spam e-mail may change significantly over time; rule-based classifiers are static (unless the user is constantly willing to make changes to the rules). Accordingly, mass e-mail senders routinely modify content of their messages in a continual attempt to prevent ("outwit") recipients from initially recognizing these messages as spam and then discarding those messages without fully reading them. Thus, unless a recipient is willing to continually construct new rules or update existing rules to track changes to spam (as that recipient perceives such changes), then, over time, a rule-based classifier becomes increasingly inaccurate at distinguishing spam from desired (non-spam) e-mail for that recipient, thereby further diminishing utility of the classifier and frustrating the user/recipient.

[0010] Alternatively, a user might consider employing a method for learning rules (as in the Cohen publication) from their existing spam in order to adapt, over time, to changes in an incoming e-mail stream. Here, the problems of a rule-based approach are more clearly highlighted. Rules are based on logical expressions; hence, as noted above, rules simply yield yes/no distinctions regarding the classification for a given e-mail message. Problematically, such rules provide no level of confidence for their predictions.

Inasmuch as users may have various tolerances as to how aggressive they would want to filter their e-mail to remove spam, then, in an application such as detecting spam, rule-based classification would become rather problematic. For example, a conservative user may require that the system be very confident that a message is spam before discarding it, whereas another user may not be so cautious. Such varying degrees of user precaution cannot be easily incorporated into a rule-based system such as that described in the Cohen publication.

## SUMMARY OF THE INVENTION

[0011] The following presents a simplified summary of the invention in order to provide a basic understanding of some aspects of the invention. This summary is not an extensive overview of the invention. It is not intended to identify key/critical elements of the invention or to delineate the scope of the invention. Its sole purpose is to present some concepts of the invention in a simplified form as a prelude to the more detailed description that is presented later.

[0012] The present invention provides for a system for detection of unsolicited messages (e.g., e-mail). The

system includes an e-mail component and a challenge component. The system can receive message(s) and associated probabilities that the message(s) are spam. Based, at least in part, upon the associated probability the system can send a challenge to a sender of a message. The e-mail component can store message(s) and associated probabilities that the messages are spam. In one example, e-mail message(s) are stored with different attributes, such as folder name, based on associated probabilities that the email message(s) are spam. In another example, e-mail message(s) having associated probabilities less than or equal to a first threshold are stored in a legitimate e-mail folder while e-mail message(s) having associated probabilities greater than the first threshold are stored in a spam folder. In yet another implementation of the invention, e-mail message(s) having associated probabilities less than or equal to a first threshold are stored in a legitimate e-mail folder, e-mail message(s) having associated probabilities greater than the first threshold, but less than or equal to a second threshold are stored in a questionable spam folder. Those e-mail message(s) having associated probabilities greater than the second threshold are stored in a spam folder. It is to be appreciated that the first threshold and/or the second threshold can be fixed, based on user preference(s) and/or adaptive (e.g., based, at least in part, upon available computational resources).

**[0013]** It will be appreciated that numbers other than probabilities, such as the score from a Support Vector Machine, a neural network, etc. can serve the same purpose as probabilities - in general, the numeric output of any machine learning algorithm can be used in place of a probability in accordance with an aspect of the present invention. Similarly, some machine learning algorithms, such as decision trees, output categorical information, and this too can be used in place of a probability combined with a threshold.

**[0014]** The challenge component can send a challenge to a sender of an e-mail message having an associated probability greater than a first threshold. For example, the challenge can be based, at least in part, upon a code embedded within the challenge (e.g., alphanumeric code). In responding to the challenge, the sender of the e-mail can reply with the code. In one example, the sender's system can be adapted to automatically retrieve the embedded code and respond to the challenge. Alternatively and/or additionally, the sender can be prompted to respond to the challenge (e.g., manually).

The use of a challenge based on an embedded code can increase the bandwidth and/or computational load of sender(s) of spam, thus, serving as a deterrent to sending of spam. It is to be appreciated that the challenge can be any of a variety of suitable types (e.g., computational challenge, a human challenge and/or a micropayment request). The challenge can be fixed and/or variable. For example, with an increased associated probability, the challenge component can send a more

difficult challenge or one that requires a greater micropayment.

**[0015]** The challenge component can modify the associated probability that the e-mail message is spam based, at least in part, upon a response to the challenge.

For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component can decrease the associated probability that the e-mail message is spam. In one example, the e-mail message is moved from a spam folder to a legitimate e-mail folder. In another implementation, the e-mail message is moved from a questionable spam folder to a legitimate e-mail folder. Upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component can increase the associated probability that the e-mail message is spam. For example, the e-mail message can be moved from a questionable spam folder to a spam folder.

**[0016]** Another aspect of the present invention provides for the system to further include a mail classifier. The mail classifier receives e-mail message(s), determines the associated probability that the e-mail message is spam and stores the e-mail message(s) and associated probabilities in the e-mail component. Accordingly, the mail classifier analyzes message content for a given recipient and distinguishes, based on that content and for that recipient, between spam and legitimate (non-spam) messages and so classifies each incoming e-mail message for that recipient.

**[0017]** Additionally and/or alternatively, e-mail message(s) can be marked with an indication of likelihood (probability) that the message is spam; message(s) assigned intermediate probabilities of spam can be moved, based on that likelihood, to questionable spam folder(s). Based, at least in part, upon information provided by the mail classifier, the challenge component can send a challenge to a sender of an e-mail message having an associated probability greater than a first threshold.

**[0018]** Yet another aspect of the present invention provides for the system to further include spam folder(s) and legitimate e-mail folder(s). The mail classifier determines the associated probability that an e-mail message is spam and stores the e-mail message in the spam folder(s) or the legitimate e-mail folder(s) (e.g., based on a first threshold). Incoming e-mail message(s) are applied to an input of the mail classifier, which, in turn, probabilistically classifies each of these messages as either legitimate or spam. Based on its classification, the message is routed to either of the spam folder(s) or the legitimate e-mail folder(s). Thereafter, the challenge component can send a challenge to a sender of an e-mail message stored in the spam folder(s) (e.g., having an associated probability greater than the first threshold). Based, at least in part, upon a response to the challenge, the challenge component can move the e-mail

message from the spam folder(s) to the legitimate e-mail folder(s). For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component can move the e-mail message from the spam folder(s) to the legitimate e-mail folder(s). Furthermore, upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component can delete the e-mail message from the spam folder(s) and/or change attribute(s) of the e-mail message stored in the spam folder(s).

**[0019]** Another aspect of the present invention provides for a system to further include a legitimate e-mail sender(s) store and/or a spam sender(s) store. The legitimate e-mail sender(s) store stores information (e.g., e-mail address) associated with sender(s) of legitimate e-mail. E-mail message(s) from sender(s) identified in the legitimate e-mail sender(s) store are generally not challenged by the challenge component. Information (e.g., e-mail address(es)) can be stored in the legitimate e-mail sender(s) store based on user selection (e.g., "do not challenge" particular sender command), a user's address book, address(es) to which a user has sent at least a specified number of e-mail messages and/or by the challenge component. The legitimate e-mail sender(s) store can further store a confidence level associated with a sender of legitimate e-mail. E-mail message(s) having associated probabilities less than or equal to the associated confidence level of the sender are not challenged by the challenge component while those e-mail message(s) having associated probabilities greater than the associated confidence level are challenged by the challenge component. The spam sender(s) store stores information (e.g., e-mail address) associated with a sender of spam. Information can be stored in the spam sender(s) store by a user and/or by the challenge component.

**[0020]** To the accomplishment of the foregoing and related ends, certain illustrative aspects of the invention are described herein in connection with the following description and the annexed drawings. These aspects are indicative, however, of but a few of the various ways in which the principles of the invention may be employed and the present invention is intended to include all such aspects and their equivalents. Other advantages and novel features of the invention may become apparent from the following detailed description of the invention when considered in conjunction with the drawings.

#### BRIEF DESCRIPTION OF THE DRAWINGS

##### **[0021]**

Fig. 1 is a block diagram of a system for detection of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 2 is a block diagram of a system for detection

of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 3 is a block diagram of a system for detection of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 4 is a block diagram of a system for detection of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 5 is a block diagram of a system for detection of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 6 is a block diagram of a system for detection of unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 7 is a block diagram of a system for responding to a challenge in accordance with an aspect of the present invention.

Fig. 8 is a flow chart illustrating a method for detecting unsolicited e-mail in accordance with an aspect of the present invention.

Fig. 9 is a flow chart further illustrating the method of Fig. 8.

Fig. 10 is a flow chart illustrating a method for responding to a challenge in accordance with an aspect of the present invention.

Fig. 11 is a flow chart illustrating a method for responding to challenges in accordance with an aspect of the present invention.

Fig. 12 is an exemplary user interface for responding to a plurality of challenges in accordance with an aspect of the present invention.

Fig. 13 illustrates an example operating environment in which the present invention may function.

#### DETAILED DESCRIPTION OF THE INVENTION

**[0022]** The present invention is now described with reference to the drawings, wherein like reference numerals are used to refer to like elements throughout. In the following description, for purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present invention. It may be evident, however, that the present invention may be practiced without these specific details. In other instances, well-known structures and devices are shown in block diagram form in order to facilitate describing the present invention.

**[0023]** As used in this application, the term "computer component" is intended to refer to a computer-related entity, either hardware, a combination of hardware and software, software, or software in execution. For example, a computer component may be, but is not limited to being, a process running on a processor, a processor, an object, an executable, a thread of execution, a program, and/or a computer. By way of illustration, both an application running on a server and the server can be a computer component. One or more computer components may reside within a process and/or thread of ex-

ecution and a component may be localized on one computer and/or distributed between two or more computers.

[0024] Referring to Fig. 1, a system 100 for detection of unsolicited messages (e.g., e-mail) in accordance with an aspect of the present invention is illustrated. The system 100 includes an e-mail component 110 and a challenge component 120. The system 100 can receive e-mail message(s) and associated probabilities that the e-mail message(s) are spam. Based, at least in part, upon the associated probability the system 100 can send a challenge to a sender of an e-mail message.

[0025] The e-mail component 110 receives and/or stores e-mail message(s) receives and/or computes associated probabilities that the e-mail messages are spam. For example, the e-mail component 110 can store information based, at least in part, upon information received from a mail classifier (not shown). In one example, e-mail message(s) are stored in the e-mail component 110 based on associated probabilities that the email message(s) are spam. In another example, the e-mail component 110 receives e-mail message(s) and computes associated probabilities that the e-mail message(s) are spam.

[0026] The challenge component 120 can send a challenge to a sender of an e-mail message having an associated probability greater than a first threshold. For example, the challenge can be based, at least in part, upon a code embedded within the challenge (e.g., alphanumeric code). In responding to the challenge, the sender of the e-mail can reply with the code. In one example, the sender's system (not shown) can be adapted to automatically retrieve the embedded code and respond to the challenge. Alternatively and/or additionally, the sender can be prompted to respond to the challenge (e.g., manually). The use of a challenge based on an embedded code can increase the bandwidth and/or computational load of sender(s) of spam, thus, serving as a deterrent to the sending of spam.

[0027] Additionally and/or alternatively the challenge can be a computational challenge, a human challenge and/or a micropayment request. These challenges and responses to these challenges are discussed more fully below. Further, the challenge can be fixed and/or variable. For example, with an increased associated probability, the challenge component 120 can send a more difficult challenge or one that requires a greater micropayment.

[0028] For example, a micropayment request can optionally utilize one-time-use spam certificates. A system 100 can put a "hold" on a received spam certificate. When a user of the system 100 reads the message and marks it as spam, the spam certificate is invalidated - sender unable to use spam certificate any further. If the message is not marked as spam, the hold is released thus allowing the sender to reuse the spam certificate (e.g., sender of message not charged money). In an alternate implementation, the spam certificate is always

invalidated at receipt, regardless of whether the message was marked as spam or not.

[0029] With regard to a computational challenge, in one implementation a challenge sender (message receiver) can determine what the computational challenge should be. However, in another implementation, the challenge is uniquely determined by some combination of the message content, the time of receipt or sending of the message, the message sender, and, importantly, the message recipient. For example, the computational challenge may be based on a one-way hash of these quantities. If the challenge sender (message recipient) is allowed to choose the challenge, than a spammer might be able to use the following technique. He subscribes to mailing lists or otherwise generates mail from users. Thus, responders send messages back to the spammer to which the spammer responds with a computational challenge of his choice. In particular, the spammer can choose challenges that legitimate users have previously sent to the spammer in response to spam! Some percentage of the recipients of the spammer's challenges solve the challenges, thus allowing the spammer to then answer the challenges sent to the spammer. In one implementation, the computational challenge is based on a one-way hash of the message (including time and recipient stamps), making it virtually impossible for sender or receiver to determine the challenge, but making it possible for each to verify that a challenge serves its intended purpose.

[0030] The challenge component 120 can modify the associated probability that the e-mail message is spam based, at least in part, upon a response to the challenge. For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component 120 can decrease the associated probability that the e-mail message is spam. In one example, the e-mail message is moved from a spam folder to a legitimate e-mail folder. In another example, the e-mail message is moved from a questionable spam folder to a legitimate e-mail folder. Moreover, upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component 120 can increase the associated probability that the e-mail message is spam.

[0031] In one implementation, a user is given a choice of challenges. For example, the choice of challenges can be based upon a filter.

[0032] Further, instead of storing the e-mail message, the system 100 can "bounce" the message, thus, necessitating the sender to resend the message along with the response to the challenge.

[0033] While Fig. 1 is a block diagram illustrating components for the system 100, it is to be appreciated that the challenge component 120 can be implemented as one or more computer components, as that term is defined herein. Thus, it is to be appreciated that computer executable components operable to implement the sys-

tem 100 and/or the challenge component 120 can be stored on computer readable media including, but not limited to, an ASIC (application specific integrated circuit), CD (compact disc), DVD (digital video disk), ROM (read only memory), floppy disk, hard disk, EEPROM (electrically erasable programmable read only memory) and memory stick in accordance with the present invention.

[0034] Turning to Fig. 2, a system 200 for detection of unsolicited e-mail in accordance with an aspect of the present invention is illustrated. The system 200 includes an e-mail component 110, a challenge component 120 and a mail classifier 130. An exemplary mail classifier 130 is set forth in greater detail in copending U.S. Patent Application entitled A TECHNIQUE WHICH UTILIZES A PROBABILISTIC CLASSIFIER TO DETECT "JUNK" E-MAIL, having serial no. 09/102,837 the entirety of which is hereby incorporated by reference. In one example, the mail classifier 130 receives e-mail message(s), determines the associated probability that the e-mail message is spam and stores the e-mail message(s) and associated probabilities in the e-mail component 110. The mail classifier 130 analyzes message content for a given recipient and distinguishes, based on that content and for that recipient, between spam and legitimate (non-spam) messages and so classifies each incoming e-mail message for that recipient.

[0035] In another example, each incoming e-mail message (in a message stream) is first analyzed to assess which one(s) of a set of predefined features, particularly characteristic of spam, the message contains. These features (e.g., the "feature set") include both simple-word-based features and handcrafted features, the latter including, for example, special multi-word phrases and various features in e-mail messages such as non-word distinctions. Generally speaking, these non-word distinctions collectively relate to, for example, formatting, authoring, delivery and/or communication attributes that, when present in a message, tend to be indicative of spam -- they are domain-specific characteristics of spam. Illustratively, formatting attributes may include whether a predefined word in the text of a message is capitalized, or whether that text contains a series of predefined punctuation marks. Delivery attributes may illustratively include whether a message contains an address of a single recipient or addresses of a plurality of recipients, or a time at which that message was transmitted (mail sent in the middle of the night is more likely to be spam). Authoring attributes may include, for example, whether a message comes from a particular e-mail address. Communication attributes can illustratively include whether a message has an attachment (a spam message rarely has an attachment), or whether the message was sent by a sender having a particular domain type (most spam appears to originate from ".com" or ".net" domain types). Handcrafted features can also include tokens or phrases known to be, for example, abusive, pornographic or insulting; or certain punc-

uation marks or groupings, such as repeated exclamation points or numbers, that are each likely to appear in spam. The specific handcrafted features are typically determined through human judgment alone or combined with an empirical analysis of distinguishing attributes of spam messages.

[0036] A feature vector, with one element for each feature in the set, is produced for each incoming e-mail message. That element simply stores a binary value specifying whether the corresponding feature is present or not in that message. The vector can be stored in a sparse format (e.g., a list of the positive features only). The contents of the vector are applied as input to a probabilistic classifier, preferably a modified support vector machine (SVM) classifier, which, based on the features that are present or absent from the message, generates a probabilistic measure as to whether that message is spam or not. This measure is then compared against a preset threshold value. If, for any message, its associated probabilistic measure equals or exceeds the threshold, then this message is classified as spam (e.g., stored in a spam folder). Alternatively, if the probabilistic measure for this message is less than the threshold, then the message is classified as legitimate (e.g., stored in a legitimate mail folder). The classification of each message can also be stored as a separate field in the vector for that message. The contents of the legitimate mail folder can then be displayed by a client e-mail program (not shown) for user selection and review. The contents of the spam folder will only be displayed by the client e-mail program upon a specific user request.

[0037] Furthermore, the mail classifier 130 can be trained using a set of M e-mail messages (e.g., a "training set", where M is an integer) that have each been manually classified as either legitimate or spam. In particular, each of these messages is analyzed to determine from a relatively large universe of n possible features (referred to herein as a "feature space"), including both simple-word-based and handcrafted features, just those particular N features (where n and N are both integers,  $n > N$ ) that are to comprise the feature set for use during subsequent classification. Specifically, a matrix (typically sparse) containing the results for all n features for the training set is reduced in size through application of Zipf's Law and mutual information, both as discussed in detail *infra* to the extent necessary, to yield a reduced N-by-m feature matrix. The resulting N features form the feature set that will be used during subsequent classification. This matrix and the known classifications for each message in the training set are then collectively applied to the mail classifier 130 for training thereof.

[0038] Furthermore, should a recipient manually move a message from one folder to another and hence reclassify it, such as from being legitimate into spam, the contents of either or both folders can be fed back as a new training set to re-train and hence update the classifier. Such re-training can occur as a result of each mes-

sage reclassification; automatically after a certain number of messages have been reclassified; after a given usage interval (e.g., several weeks or months) has elapsed; or upon user request. In this manner, the behavior of the classifier can advantageously track changing subjective perceptions and preferences of its particular user. Alternatively, e-mail messages may be classified into multiple categories (subclasses) of spam (e.g., commercial spam, pornographic spam and so forth). In addition, messages may be classified into categories corresponding to different degrees of spam (e.g., "certain spam", "questionable spam", and "non-spam").

**[0039]** Based, at least in part, upon information provided by the mail classifier 130, the challenge component 120 can send a challenge to a sender of an e-mail message having an associated probability greater than a first threshold. For example, the challenge can be based, at least in part, upon a code embedded within the challenge (e.g., alphanumeric code). In responding to the challenge, the sender of the e-mail can reply with the code. The sender's system (not shown) can be adapted to automatically retrieve the embedded code and respond to the challenge. Alternatively and/or additionally, the sender can be prompted to respond to the challenge (e.g., manually). The use of a challenge based on an embedded code can increase the bandwidth and/or computational load of sender(s) of spam, thus, serving as a deterrent to the sending of spam. It is to be appreciated that any type of challenge (e.g., a computational challenge, a human challenge, a micropayment request) suitable for carrying out the present invention can be employed and all such types of challenges are intended to fall within the scope of the hereto appended claims.

**[0040]** The challenge component 120 can modify the associated probability that an e-mail message is spam based, at least in part, upon a response to the challenge. For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component 120 can decrease the associated probability that the e-mail message is spam.

**[0041]** Upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component 120 can increase the associated probability that the e-mail message is spam. It is to be appreciated that the mail classifier 130 can be a computer component as that term is defined herein.

**[0042]** Referring next to Fig. 3, a system 300 for detection of unsolicited e-mail in accordance with an aspect of the present invention is illustrated. The system 300 includes a mail classifier 310, a challenge component 320, spam folder(s) 330 and legitimate e-mail folder(s) 340. In one implementation, the spam folder(s) 330 and/or the legitimate e-mail folder(s) 340 can be virtual, that is, storing information associated with e-mail message(s) (e.g., link to e-mail message(s)) with the e-mail

message(s) stored elsewhere. Or, in another implementation, rather than folders, an attribute of the message, can simply be set.

**[0043]** As discussed *supra*, the mail classifier 310 determines the associated probability that an e-mail message is spam and stores the e-mail message in the spam folder(s) 330 or the legitimate e-mail folder(s) 340 (e.g., based on a first threshold). Incoming e-mail message(s) are applied to an input of the mail classifier 310, which, in turn, probabilistically classifies each of these messages as either legitimate or spam. Based on its classification, the e-mail message is routed to either of the spam folder(s) 330 or the legitimate e-mail folder(s) 340. Thus, e-mail message(s) having associated probabilities less than or equal to a first threshold are stored in a legitimate e-mail folder(s) 340 while e-mail message(s) having associated probabilities greater than the first threshold are stored in a spam folder(s) 330. The first threshold can be fixed, based on user preference(s) and/or adaptive (e.g., based, at least in part, upon available computational resources).

**[0044]** Thereafter, the challenge component 320 can send a challenge to a sender of an e-mail message stored in the spam folder(s) (e.g., having an associated probability greater than the first threshold). For example, the challenge can be based, at least in part, upon a code embedded within the challenge, a computational challenge, a human challenge and/or a micropayment request. Based, at least in part, upon a response to the challenge, the challenge component 320 can move the e-mail message from the spam folder(s) 330 to the legitimate e-mail folder(s) 340. For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component 320 can move the e-mail message from the spam folder(s) 330 to the legitimate e-mail folder(s) 340.

**[0045]** Upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component 320 can delete the e-mail message from the spam folder(s) 330 and/or change attribute(s) of the e-mail message stored in the spam folder(s) 330. For example, display attribute(s) (e.g., color) of the e-mail message can be changed to bring to a user's attention the increased likelihood of the e-mail message being spam.

**[0046]** Next, turning to Fig. 4, a system 400 for detection of unsolicited e-mail in accordance with an aspect of the present invention is illustrated. The system 400 includes a mail classifier 310, a challenge component 320, spam folder(s) 330 and legitimate e-mail folder(s) 340. The system 400 further includes a legitimate e-mail sender(s) store 350 and/or a spam sender(s) store 360. The legitimate e-mail sender(s) store 350 stores information (e.g., e-mail address) associated with sender(s) of legitimate e-mail. E-mail message(s) from sender(s) identified in the legitimate e-mail sender(s) store 350 are generally not challenged by the challenge component

320. Accordingly, in one example, e-mail message(s) stored in the spam folder(s) 330 by the mail classifier 310 are moved to the legitimate mail folder(s) 340 if the sender of the e-mail message is stored in the legitimate e-mail sender(s) store 350.

[0047] Information (e.g., e-mail address(es)) can be stored in the legitimate e-mail sender(s) store 350 based on user selection (e.g., "do not challenge" particular sender command), a user's address book, address(es) to which a user has sent at least a specified number of e-mail messages and/or by the challenge component 320. For example, once a sender of an e-mail message has responded correctly to a challenge, the challenge component 320 can store information associated with the sender (e.g., e-mail address) in the legitimate e-mail sender(s) store 350.

[0048] The legitimate e-mail sender(s) store 350 can further retain a confidence level associated with a sender of legitimate e-mail. E-mail message(s) having associated probabilities less than or equal to the associated confidence level of the sender are not challenged by the challenge component 320 while those e-mail message(s) having associated probabilities greater than the associated confidence level are challenged by the challenge component 320. For example, the confidence level can be based, at least in part, upon the highest associated probability challenge to which the sender has responded.

[0049] In one implementation, a sender can be removed from the legitimate e-mail sender(s) store 350 based, at least in part, upon a user's action (e.g., e-mail message from the sender deleted as spam). In accordance with another aspect, sender(s) are added to the legitimate e-mail sender(s) store 350 after a user has sent one e-mail message to the sender - this can be useful for mailing list(s).

[0050] The spam sender(s) store 360 stores information (e.g., e-mail address) associated with a sender of spam. Information can be stored in the spam sender(s) store 360 by a user and/or by the challenge component 320. For example, once a user has deleted a particular e-mail message as spam, information associated with the sender of the e-mail message can be stored in the spam sender(s) store 360. In another example, information associated with a sender of an e-mail message that incorrectly responded to a challenge and/or failed to respond to the challenge can be stored in the spam sender(s) store 360.

[0051] Fig. 5 illustrates a system 500 for detection of unsolicited e-mail in accordance with an aspect of the present invention is illustrated. The system 500 includes a mail classifier 510, a challenge component 520, spam folder(s) 530, questionable spam folder(s) 540 and legitimate e-mail folder(s) 550. As discussed above, the mail classifier 510 determines the associated probability that an e-mail message is spam and stores the e-mail message in the spam folder(s) 530, the questionable spam folder(s) 540 or the legitimate e-mail folder(s) 550.

Incoming e-mail message(s) are applied to an input of the mail classifier 510, which, in turn, probabilistically classifies each of these messages as either legitimate, questionable spam or spam. Based on its classification, each message is routed to one of the spam folder(s) 530, the questionable spam folder(s) 540 or the legitimate e-mail folder(s) 550.

[0052] E-mail message(s) having associated probabilities less than or equal to a first threshold are in legitimate e-mail folder(s) 550. E-mail message(s) having associated probabilities greater than the first threshold, but less than or equal to a second threshold are stored in questionable spam folder(s) 540. Further, e-mail message(s) having associated probabilities greater than the second threshold are stored in spam folder(s) 530. It is to be appreciated that the first threshold and/or the second threshold can be fixed, based on user preference (s) and/or adaptive (e.g., based, at least in part, upon available computational resources). Thereafter, the challenge component 520 can send a challenge to a sender of an e-mail message stored in the questionable spam folder(s) 540. For example, the challenge can be based, at least in part, upon a code embedded within the challenge, a computational challenge, a human challenge and/or a micropayment request.

[0053] Based, at least in part, upon a response to the challenge or lack thereof, the challenge component 520 can move the e-mail message from the questionable spam folder(s) 540 to the legitimate e-mail folder(s) 550 or the spam folder(s) 530. For example, upon receipt of an appropriate (e.g., correct) response to the challenge, the challenge component 520 can move the e-mail message from the questionable spam folder(s) 540 to the legitimate e-mail folder(s) 550.

[0054] Further, upon receipt of an inappropriate (e.g., incorrect) response to the challenge and/or failure to receive a response to the challenge in a particular time period (e.g., 4 hours), the challenge component 520 can move the e-mail message from the questionable spam folder(s) 540 to the spam folder(s) 530.

[0055] Referring next to Fig. 6, a system 600 for detection of unsolicited e-mail in accordance with an aspect of the present invention is illustrated. The system 600 includes a mail classifier 510, a challenge component 520, spam folder(s) 530, questionable spam folder(s) 540 and legitimate e-mail folder(s) 550. The system 600 further includes a legitimate e-mail sender(s) store 560 and/or a spam sender(s) store 570.

[0056] The legitimate e-mail sender(s) store 560 stores information (e.g., e-mail address) associated with sender(s) of legitimate e-mail. E-mail message(s) from entities identified in the legitimate e-mail sender(s) store 560 are generally not challenged by the challenge component 520. Accordingly, in one example, e-mail message(s) stored in the spam folder(s) 530 or the questionable spam folder(s) 540 by the mail classifier 510 are moved to the legitimate mail folder(s) 550 if the sender of the e-mail message is stored in the legitimate

e-mail sender(s) store 560.

[0057] Information (e.g., e-mail address(es)) can be stored in the legitimate e-mail sender(s) store 660 based on user selection (e.g., "do not challenge" particular sender command), a user's address book, address(es) to which a user has sent at least a specified number of e-mail messages and/or by the challenge component 520. For example, once a sender of an e-mail message has responded correctly to a challenge, the challenge component 520 can store information associated with the sender (e.g., e-mail address) in the legitimate e-mail sender(s) store 560.

[0058] The legitimate e-mail sender(s) store 560 can further store a confidence level associated with a sender of legitimate e-mail. For example, e-mail message(s) having associated probabilities less than or equal to the associated confidence level of the sender are not challenged by the challenge component 520 while those e-mail message(s) having associated probabilities greater than the associated confidence level are challenged by the challenge component 520. For example, the confidence level can be based, at least in part, upon the highest associated probability challenge to which the sender has responded.

[0059] In one example, a sender can be removed from the legitimate e-mail sender(s) store 560 based, at least in part, upon a user's action (e.g., e-mail message from the sender deleted as spam). In another example, sender(s) are added to the legitimate e-mail sender(s) store 560 after a user has sent one e-mail message to the sender.

[0060] The spam sender(s) store 570 stores information (e.g., e-mail address) associated with a sender of spam. Information can be stored in the spam sender(s) store 570 by a user and/or by the challenge component 520. For example, once a user has deleted a particular e-mail message as spam, information associated with the sender of the e-mail message can be stored in the spam sender(s) store 570. In another example, information associated with a sender of an e-mail message that incorrectly responded to a challenge and/or failed to respond to the challenge can be stored in the spam sender(s) store 570.

[0061] In one example, a unique-ID can be exchanged during the challenge process (e.g., to reduce the likelihood that a spammer can send spam using an address of a good sender). Further, sender(s) can use message signing. Unsigned message(s) from sender(s) stored in the legitimate e-mail sender(s) store 560 who usually sign their message(s) are subjected to the usual processing and potential challenging.

[0062] In another example, higher volume sender(s) of e-mail customize their "from" address (e.g., a unique "from" address for a recipient). For example, the "from" address can be based on a global secret key known to the sender and hashed with the recipient's e-mail address. Alternatively, a random number can be generated and stored for a recipient.

[0063] In yet a third example, a "per recipient ID" (PRID) is included in e-mail message(s). The PRID appends sender unique information in a special message header field. It is to be appreciated that the PRID does not have to be set on a per-sender basis. Thus, as mail is forwarded around an organization, inclusion on the legitimate e-mail sender(s) store 560 can be inherited. The PRID can be a public key for use with a public key signature system (e.g., OpenPGP or S/MIME).

[0064] Additionally, sender(s) of e-mail message(s) can include requests for challenge(s) (e.g., to facilitate scheduling of receipt of challenge(s)). For example, an e-mail message(s) can include a "CHALLENGE\_ME\_NOW: TRUE" header. This can cause a system 600 to automatically send a challenge and when a correct response is received to include the sender in the legitimate e-mail sender(s) store 560.

[0065] The challenge component 520 can be adapted to detect e-mail message(s) received from mailing list(s) (e.g., moderated mailing list(s) and/or unmoderated mailing list(s)). For example, a header line such as "Precedence: list" or "Precedence: bulk" can be included in e-mail message(s) received from a mailing list. In another example, the challenge component 520 detects that an e-mail message is spam based, at least in part upon, detection of a "sender" line being different from a "from" line. E-mail message header(s) typically contains two different from lines: one "from" line at the top (e.g., inserted by the from command used by SMTP), and a "from:" header field (e.g., the one that is usually displayed to the user.) For mailing lists, these may differ.

[0066] In one example, the challenge component 520 can detect e-mail message(s) from mailing list(s) and give a user the opportunity to include the mailing list(s) in the legitimate e-mail sender(s) store 560. The challenge component 520 can additionally include a level of confidence associated with the mailing list(s).

[0067] An issue to be addressed with regard to mailing list(s) is to reduce the likelihood that spam-like message(s) received from a mailing list will create a mail storm of challenges to the mailing list. This issue differs for the different list types. There are 8 situations, although many of them share the same solution. In particular, a mailing list can be moderated or unmoderated and additionally can have different level(s) of ability to respond to challenges. This creates 8 types.

[0068] Many moderated mailing list(s) include an "approved-by" header. For example, for moderated mailing list(s), it can be assumed that either all messages are good, or all are spam. For unmoderated lists, it can be assumed that some spam will be sent to the mailing list. Thus, for an unmoderated mailing list, the challenge component 520 can allow a user to set a threshold determining whether spam-like messages should be shown, or simply put in the spam folder(s) 530.

[0069] For example, an e-mail message from a mailing list has been detected, a user is given the user the opportunity to determine the level of confidence associ-

ated with the mailing list. A concern is sending too many challenges to mailing lists, especially those that do not have the ability to automatically respond to challenges. For moderated mailing list(s), for example, a user can be prompted to include the mailing list in the legitimate e-mail sender(s) store 560. In another example, the mailing list can respond to a challenge from the challenge component 520 and be included in the legitimate e-mail sender(s) store 560. In yet a third example, upon subscription to the mailing list, the mailing list prompts the user to include the mailing list in the user's legitimate e-mail sender(s) store 560.

**[0070]** For unmoderated mailing list(s), for example, a user can be prompted to set a threshold for the mailing list. E-mail message(s) having a probability of being spam above the threshold is moved to the spam folder (s) 530 and/or deleted. In another example, the mailing list can respond to a challenge from the challenge component 520 and be included in the legitimate e-mail sender(s) store 560. In yet a third example, upon subscription to the mailing list, the mailing list prompts the user to include the mailing list in the user's legitimate e-mail sender(s) store 560.

**[0071]** The challenge component 520 can take into account mailing list(s) that do not have the ability to automatically respond to challenges. In particular, for moderated mailing lists, the challenge component 520 can include the mailing list in the legitimate e-mail sender(s) store 560. For unmoderated mailing lists, the challenge component 520 can facilitate setting a threshold for the mailing list: messages above the threshold are challenged while messages below the threshold are let through

**[0072]** Inclusion in the legitimate e-mail sender(s) store 560 can occur at an appropriate time. For mailing lists, it is likely that the user will not send mail TO the list. However, it is undesirable to include the mailing list in the legitimate e-mail sender(s) store 560 based on small amounts of mail received FROM the list. Otherwise a spammer could masquerade as a mailing list, send a small amount of messages (none of which are deleted as spam) and then send spam freely. In one implementation, the first time that mail from a mailing list arrives, and is not detected as spam, the user is prompted to add the mailing list to the legitimate e-mail sender (s) store 560, with an associated threshold. Since most mailing lists include a welcome message, if some welcome messages are included in training data, the welcome message is unlikely to be marked as spam.

**[0073]** If, however, the first messages that arrive are substantially all spam-like, then the messages should be included in the spam folder(s) 530. In particular, it is not desirable to let someone masquerade as a mailing list, and send spam. Thus, until the mail listing is included in the legitimate e-mail sender(s) store 560, the challenge component 520 can send challenge(s) to the mailing list as described *supra*. If the messages are spam-like but legitimate, the user may or may not receive

them, depending on how the challenges are handled. If the challenges are not answered, they will not get through. Thus, it should be difficult to get spam through. Eventually, the mailing list will send a non-spam like message, and the user will be prompted to establish a policy for the mailing list.

**[0074]** It is to be appreciated that mailing list(s) may have a From address such that mail sent to that From address is sent to the entire list. If a list appears to be of that type, it is undesirable to send challenges to it as they might be received by substantially all readers of the mailing list. Apparent spam from such a mailing list before the mailing list has been included in the legitimate e-mail sender(s) store 560 can simply be ignored.

The definition of inclusion in the legitimate e-mail sender (s) store 560 can be modified for mailing list(s). Given that the From line on a mailing list, even a moderated one is different for each sender, inclusion in the legitimate e-mail sender(s) store 560 can be based on other part(s) of the header. Often, the To line on a mailing list is the mailing list name (so that reply-all goes to the whole list.). Thus, for mailing lists, inclusion in the legitimate e-mail sender(s) store 560 can be based, at least in part, on the to-line. This can be in addition to from-line listing (e.g., if the sender of the mailing list is in the legitimate e-mail sender(s) store 560 that also should be sufficient). It is to be appreciated that other header lines, for mailing lists, such as sent-by, that can additionally and/or alternatively be included in the legitimate e-mail sender(s) store 560.

**[0075]** In order to determine validity of e-mail address (es), spammer(s) rely on "bouncing". Many conventional e-mail servers bounce e-mail back to it's sender if it is addressed to an invalid address. Thus, for e-mail servers those e-mail servers, the indicia of validity of an e-mail address increases if an e-mail message is not bounced. Accordingly, spammers can send more spam messages to the unbounced addresses.

**[0076]** For those e-mail servers which bounce e-mail, challenges of the present invention do not provide any additional information to the spammer (e.g., lack of bounce is an indication of validity of the address). Further, the e-mail server can itself send challenges via a system for detection of unsolicited e-mail for "semi-live" address(es) (e.g., valid but unmonitored address).

**[0077]** With regard to e-mail servers which do not bounce e-mail to invalid addresses, again the e-mail server can itself send challenges via a system for detection of unsolicited e-mail, for example, to have behavior of invalid address(es) be similar to the behavior of valid address(es). Further, in one implementation, a randomization factor is added to the probability that an e-mail is spam by the server system (e.g., to prevent attempts to circumvent adaptive spam filters).

**[0078]** Next, turning to Fig. 7, a system 700 for responding to a challenge in accordance with an aspect of the present invention is illustrated. The system 700 includes a challenge receiver component 710, a chal-

challenge processor component 720 and a challenge response component 730:

**[0079]** The challenge receiver component 710 receives a challenge (e.g., to a previously sent e-mail). For example the challenge can be based, at least in part, upon a code embedded within the challenge, a computational challenge, a human challenge and/or a micropayment request.

**[0080]** In one example, the challenge receiver component 710 determines which of a plurality of challenge modalities to forward to the challenge processor component 720 (e.g., based on available computational resources and/or user preference). In another example, the challenge receiver component 710 provides information to a user to facilitate selection of one of a plurality of challenge modalities, thus, allowing a user to select which modality, if any, the user wishes to use to respond to the challenge. For example, the challenge receiver component 710 can provide information which may be helpful to the user in selecting an appropriate response modality, such as, an amount of computational resources required to respond to a computational challenge, an amount of a micropayment and/or a balance of a micropayment account. Once a challenge modality has been selected, the challenge is forwarded to the challenge processor 720.

**[0081]** It is to be appreciated that in certain instances the user may desire to not respond to the challenge, in which case, no information is sent to the challenge processor component 720 and/or the challenge response component 730.

**[0082]** The challenge processor component 720 processes the challenge and provides an output associated with the processed challenge. For example, when the challenge includes an embedded code, the challenge processor component 720 can provide an output to the challenge response component 730 which includes the embedded code. In the instance in which the challenge includes a computational challenge, the challenge processor component 720 can facilitate generation of a solution to the computational challenge.

**[0083]** When the challenge includes a human challenge, the challenge processor component 720 can provide information to a user to facilitate solving the human challenge. In one example, the human challenge can include a problem that is relatively easy for a human to solve, and relatively hard for a computer. In one example, the human challenge includes an image of a word (e.g., GIF or JPEG). The word is partially obscured by noise. The noise makes it hard to automatically develop a computer program to read the word (or at least, to use off-the-shelf components), without making it too hard for a human to do it. In this example, the challenge processor component 720 can provide the image of the word to the user. The user then provides the word back to the challenge processor component 720. The challenge processor component 720 provides an output including the word to the challenge response component 730.

**[0084]** When the challenge includes a micropayment request, the challenge processor component 720 can facilitate providing an output to the challenge response component 730. In one example, a response to a micropayment challenge is based on a one-time use "spam certificate" which can be issued by an issuing authority.

The challenge processor component 720 can either automatically or based on user input provides a spam certificate number to the challenge response component 730. By providing the spam certificate number, the spam certificate is thereafter invalidated (e.g., one-time use).

**[0085]** In another example, a response to a micropayment challenge is based on a micropayment account. Each such response causes an amount to be removed from a micropayment account maintained, for example, by an issuing authority. The challenge processor component 720 can provide information associated with the micropayment account to the challenge response component 730.

**[0086]** The challenge response component 730 provides a response to the challenge based, at least in part, upon the output associated with the processed challenge. For example, the response to the challenge can include an embedded code, solution to a computational challenge, solution to a human challenge and/or micropayment.

**[0087]** In one implementation, for example, to reduce a likelihood of a denial-of-service attack, computational challenges are ordered by the quantity of challenges already processed for a given message. Message(s) with fewer processed challenge(s) are processed before message(s) having a greater quantity of processed challenges are processed (e.g., as computational resources are available). Thus, in the instance in which a message is sent to a mailing list, a recipient could send computational challenges in an effort to maliciously cause a denial-of-service attack. However, once one or more computational challenges are processed for that message, computational challenges of other messages having less processed challenges are given priority, thus reducing the likelihood of a denial-of-service.

**[0088]** In view of the exemplary systems shown and described above, methodologies that may be implemented in accordance with the present invention will be better appreciated with reference to the flow chart of Figs. 8, 9, 10 and 11. While, for purposes of simplicity of explanation, the methodologies are shown and described as a series of blocks, it is to be understood and appreciated that the present invention is not limited by the order of the blocks, as some blocks may, in accordance with the present invention, occur in different orders and/or concurrently with other blocks from that shown and described herein. Moreover, not all illustrated blocks may be required to implement the methodologies in accordance with the present invention.

**[0089]** The invention may be described in the general context of computer-executable instructions, such as program modules, executed by one or more compo-

nents. Generally, program modules include routines, programs, objects, data structures, *etc.* that perform particular tasks or implement particular abstract data types. Typically the functionality of the program modules may be combined or distributed as desired in various embodiments.

**[0090]** Turning to Figs. 8 and 9, a method 800 for detecting an unsolicited e-mail message in accordance with an aspect of the present invention is illustrated. At 804, an e-mail message is received. At 808, a probability that the e-mail message is spam is determined (e.g., by a mail classifier).

**[0091]** At 812, a determination is made as to whether the sender of the e-mail message is in a legitimate e-mail sender(s) store. If the determination at 812 is YES, processing continues at 816. If the determination at 812 is NO, at 820, a determination is made as to whether the sender of the e-mail message is in a spam sender (s) store. If the determination at 820 is YES, processing continues at 824. If the determination at 820 is NO, at 828, a determination is made as to whether the probability that the e-mail message is spam is greater than a first threshold. If the determination at 828 is NO, processing continues at 816. If the determination at 828 is YES, at 832, one or more challenge(s) are sent to the sender of the e-mail message.

**[0092]** At 836, a determination is made as to whether a response to the challenge(s) has been received. If the determination at 836 is NO, processing continues at 836. If the determination at 836 is YES, at 840, a determination is made as to whether the response received to the challenge is correct. If the determination at 840 is YES, processing continues at 816. If the determination at 840 is NO, processing continues at 824.

**[0093]** At 816, the e-mail message is identified as "not spam" (e.g., placed in legitimate e-mail folder(s) and/or associated probability decreased). Next, at 844, the sender of the e-mail message is added to the legitimate e-mail sender(s) store and no further processing occurs.

**[0094]** At 824, the e-mail message is identified as spam (e.g., placed in spam folder(s), deleted and/or associated probability increased). Next, at 848, the sender of the e-mail message is added to the spam sender(s) store and no further processing occurs.

**[0095]** Referring next to Fig. 10, a method 1000 for responding to a challenge in accordance with an aspect of the present invention is illustrated. At 1010, an e-mail message is sent. At 1020, a challenge is received (e.g., an embedded code, a computational challenge, a human challenge and/or a request for a micropayment). At 1030, the challenge is processed. At 1040, a response to the challenge is sent.

**[0096]** Next, turning to Fig. 11, a method 1100 for responding to challenges in accordance with an aspect of the present invention is illustrated. At 1110, e-mail message(s) are sent. At 1120, challenge(s) are received (e.g., each challenge having an embedded code, a com-

putational challenge, a human challenge and/or a request for a micropayment). At 1130, the challenge(s) to be processed are ordered based, at least in part, upon message(s) with fewer challenge(s) processed before message(s) with more challenge(s) processed (e.g., to reduce denial-of-service attacks). At 1140, the challenge is processed. At 1150, a response to the selected challenge is sent. At 1160, a determination is made as to whether there are more challenge(s) to process. If the determination at 1160 is YES, processing continues at 1130. If the determination at 1160 is NO, no further processing occurs.

**[0097]** Turning to Fig. 12, an exemplary user interface 1200 for responding to a plurality of challenges in accordance with an aspect of the present invention is illustrated. In this exemplary user interface, a user is prompted with the message:

THE E-MAIL MESSAGE YOU SENT HAS BEEN DETECTED AS POTENTIAL SPAM. UNLESS YOU CORRECTLY RESPOND TO ONE OF THE CHALLENGES IDENTIFIED BELOW, THE E-MAIL MESSAGE MAY BE IDENTIFIED AS SPAM AND/OR DELETED AS SPAM.

**[0098]** The user is presented with three options: computer computational challenge, human challenge and micropayment. Based, at least in part, upon the user's selection, the selected challenge can then be processed.

**[0099]** In order to provide additional context for various aspects of the present invention, Fig. 13 and the following discussion are intended to provide a brief, general description of a suitable operating environment 1310 in which various aspects of the present invention may be implemented. While the invention is described in the general context of computer-executable instructions, such as program modules, executed by one or more computers or other devices, those skilled in the art will recognize that the invention can also be implemented in combination with other program modules and/or as a combination of hardware and software. Generally, however, program modules include routines, programs, objects, components, data structures, *etc.* that perform particular tasks or implement particular data types. The operating environment 1310 is only one example of a suitable operating environment and is not intended to suggest any limitation as to the scope of use or functionality of the invention. Other well known computer systems, environments, and/or configurations that may be suitable for use with the invention include but are not limited to, personal computers, hand-held or laptop devices, multiprocessor systems, microprocessor-based systems, programmable consumer electronics, network PCs, minicomputers, mainframe computers, distributed computing environments that include the above systems or devices, and the like.

**[0100]** With reference to Fig. 13, an exemplary envi-

environment 1310 for implementing various aspects of the invention includes a computer 1312. The computer 1312 includes a processing unit 1314, a system memory 1316, and a system bus 1318. The system bus 1318 couples system components including, but not limited to, the system memory 1316 to the processing unit 1314. The processing unit 1314 can be any of various available processors. Dual microprocessors and other multiprocessor architectures also can be employed as the processing unit 1314.

**[0101]** The system bus 1318 can be any of several types of bus structure(s) including the memory bus or memory controller, a peripheral bus or external bus, and/or a local bus using any variety of available bus architectures including, but not limited to, 13-bit bus, Industrial Standard Architecture (ISA), Micro-Channel Architecture (MSA), Extended ISA (EISA), Intelligent Drive Electronics (IDE), VESA Local Bus (VLB), Peripheral Component Interconnect (PCI), Universal Serial Bus (USB), Advanced Graphics Port (AGP), Personal Computer Memory Card International Association bus (PCMCIA), and Small Computer Systems Interface (SCSI).

**[0102]** The system memory 1316 includes volatile memory 1320 and nonvolatile memory 1322. The basic input/output system (BIOS), containing the basic routines to transfer information between elements within the computer 1312, such as during start-up, is stored in nonvolatile memory 1322. By way of illustration, and not limitation, nonvolatile memory 1322 can include read only memory (ROM), programmable ROM (PROM), electrically programmable ROM (EPROM), electrically erasable ROM (EEPROM), or flash memory. Volatile memory 1320 includes random access memory (RAM), which acts as external cache memory. By way of illustration and not limitation, RAM is available in many forms such as synchronous RAM (SRAM), dynamic RAM (DRAM), synchronous DRAM (SDRAM), double data rate SDRAM (DDR SDRAM), enhanced SDRAM (ESDRAM), Synchlink DRAM (SLDRAM), and direct Rambus RAM (DRRAM).

**[0103]** Computer 1312 also includes removable/non-removable, volatile/nonvolatile computer storage media. Fig. 13 illustrates, for example a disk storage 1324. Disk storage 1324 includes, but is not limited to, devices like a magnetic disk drive, floppy disk drive, tape drive, Jaz drive, Zip drive, LS-100 drive, flash memory card, or memory stick. In addition, disk storage 1324 can include storage media separately or in combination with other storage media including, but not limited to, an optical disk drive such as a compact disk ROM device (CD-ROM), CD recordable drive (CD-R Drive), CD rewritable drive (CD-RW Drive) or a digital versatile disk ROM drive (DVD-ROM). To facilitate connection of the disk storage devices 1324 to the system bus 1318, a removable or non-removable interface is typically used such as interface 1326.

**[0104]** It is to be appreciated that Fig 13 describes software that acts as an intermediary between users

and the basic computer resources described in suitable operating environment 1310. Such software includes an operating system 1328. Operating system 1328, which can be stored on disk storage 1324, acts to control and allocate resources of the computer system 1312. System applications 1330 take advantage of the management of resources by operating system 1328 through program modules 1332 and program data 1334 stored either in system memory 1316 or on disk storage 1324. It is to be appreciated that the present invention can be implemented with various operating systems or combinations of operating systems.

**[0105]** A user enters commands or information into the computer 12 through input device(s) 1336. Input devices 1336 include, but are not limited to, a pointing device such as a mouse, trackball, stylus, touch pad, keyboard, microphone, joystick, game pad, satellite dish, scanner, TV tuner card, digital camera, digital video camera, web camera, and the like. These and other input devices connect to the processing unit 1314 through the system bus 1318 via interface port(s) 1338. Interface port(s) 1338 include, for example, a serial port, a parallel port, a game port, and a universal serial bus (USB). Output device(s) 1340 use some of the same type of ports as input device(s) 1336. Thus, for example, a USB port may be used to provide input to computer 1312, and to output information from computer 1312 to an output device 1340. Output adapter 1342 is provided to illustrate that there are some output devices 1340 like monitors, speakers, and printers among other output devices 1340 that require special adapters. The output adapters 1342 include, by way of illustration and not limitation, video and sound cards that provide a means of connection between the output device 1340 and the system bus 1318. It should be noted that other devices and/or systems of devices provide both input and output capabilities such as remote computer(s) 1344.

**[0106]** Computer 1312 can operate in a networked environment using logical connections to one or more remote computers, such as remote computer(s) 1344. The remote computer(s) 1344 can be a personal computer, a server, a router, a network PC, a workstation, a microprocessor based appliance, a peer device or other common network node and the like, and typically includes many or all of the elements described relative to computer 1312. For purposes of brevity, only a memory storage device 1346 is illustrated with remote computer(s) 1344. Remote computer(s) 1344 is logically connected to computer 1312 through a network interface 1348 and then physically connected via communication connection 1350. Network interface 1348 encompasses communication networks such as local-area networks (LAN) and wide-area networks (WAN). LAN technologies include Fiber Distributed Data Interface (FDDI), Copper Distributed Data Interface (CDDI), Ethernet/IEEE 1302.3, Token Ring/IEEE 1302.5 and the like. WAN technologies include, but are not limited to, point-to-point links, circuit switching networks like Integrated

Services Digital Networks (ISDN) and variations thereon, packet switching networks, and Digital Subscriber Lines (DSL).

[0107] Communication connection(s) 1350 refers to the hardware/software employed to connect the network interface 1348 to the bus 1318. While communication connection 1350 is shown for illustrative clarity inside computer 1312, it can also be external to computer 1312. The hardware/software necessary for connection to the network interface 1348 includes, for exemplary purposes only, internal and external technologies such as, modems including regular telephone grade modems, cable modems and DSL modems, ISDN adapters, and Ethernet cards.

[0108] What has been described above includes examples of the present invention. It is, of course, not possible to describe every conceivable combination of components or methodologies for purposes of describing the present invention, but one of ordinary skill in the art may recognize that many further combinations and permutations of the present invention are possible. Accordingly, the present invention is intended to embrace all such alterations, modifications and variations that fall within the spirit and scope of the appended claims. Furthermore, to the extent that the term "includes" is used in either the detailed description or the claims, such term is intended to be inclusive in a manner similar to the term "comprising" as "comprising" is interpreted when employed as a transitional word in a claim.

#### Claims

1. A system facilitating detection of unsolicited e-mail, comprising:
  - an e-mail component that receives or stores messages and receives or computes associated probabilities that the e-mail messages are spam; and,
  - a challenge component that sends a challenge to an originator of an e-mail message having an associated probability greater than a first threshold.
2. The system of claim 1, further comprising a mail classifier that receives e-mail messages and determines the associated probability that the e-mail message is spam.
3. The system of claim 1, the challenge component further modifying the associated probability that the e-mail message is spam based, at least in part, upon a response to the challenge.
4. The system of claim 1, the challenge being an embedded code.

5. The system of claim 1, the challenge being a computational challenge.
6. The system of claim 5, the computational challenge being a one-way hash of the message including time stamp and recipient stamp.
7. The system of claim 1, the challenge being a human challenge.
8. The system of claim 1, the challenge being a micro-payment request.
9. The system of claim 1, a user being given a choice of challenges, the choice of challenges being based upon a filter.
10. The system of claim 1, a difficulty of the challenge being based, at least in part, upon the associated probability that the e-mail message is spam.
11. A system facilitating detection of unsolicited messages, comprising:
  - a mail classifier that receives an incoming message and classifies the incoming message as spam or a legitimate message; and,
  - a challenge component that sends a challenge to a sender of the message if the message is classified as spam.
12. The system of claim 11, the mail classifier further storing the incoming message in a spam folder or a legitimate message folder.
13. The system of claim 12, the challenge component further moving the message from the spam folder to the legitimate message folder based, at least in part, upon a response to the challenge.
14. The system of claim 11, the challenge being an embedded code.
15. The system of claim 11, the challenge being a computational challenge.
16. The system of claim 11, the challenge being a human challenge.
17. The system of claim 11, the challenge being a micro-payment request.
18. The system of claim 11, further comprising a legitimate message sender(s) store that stores information associated with a sender of legitimate message(s).
19. The system of claim 18, the challenge component

adding information associated with the sender of the message to the legitimate message sender(s) store, if the challenge is responded to correctly.

20. The system of claim 11, further comprising a spam sender(s) store that stores information associated with a sender of spam. 5
21. A system facilitating detection of unsolicited e-mail, comprising: 10
  - a mail classifier that receives an incoming e-mail message and classifies the incoming e-mail message as spam, questionable spam or legitimate e-mail; and, 15
  - a challenge component that sends a challenge to a sender of an e-mail message classified as questionable spam.
22. The system of claim 21, the mail classifier further storing the incoming e-mail message in a spam folder, a questionable spam or a legitimate mail folder. 20
23. The system of claim 22, the challenge component further moving the e-mail message from the questionable spam folder to the spam folder or the legitimate mail folder based, at least in part, upon a response to the challenge. 25
24. The system of claim 21, the challenge being at least one of an embedded code, a computational challenge, a human challenge and a micropayment request. 30
25. The system of claim 21 further comprising a legitimate e-mail sender(s) store that stores information associated with a sender of legitimate e-mail. 35
26. The system of claim 21, further comprising a spam sender(s) store that stores information associated with a sender of spam. 40
27. The system of claim 21, the e-mail message including a per recipient ID. 45
28. The system of claim 21, the challenge component further adapted to detect whether the e-mail message is from a mailing list.
29. The system of claim 28, the challenge component further adapted to detect whether the mailing list is moderated or unmoderated. 50
30. A method for detecting unsolicited e-mail, comprising: 55

sending a challenge to a sender of an e-mail message classified as questionable spam;

receiving a response to the challenge; and, modifying the classification of the e-mail message based, at least in part, upon the response to the challenge.

31. The method of claim 30, further comprising at least one of the following acts, receiving the e-mail message;
  - classifying the e-mail message as spam, questionable spam or legitimate e-mail;
  - determining whether the sender is stored in a legitimate e-mail sender(s) store; and,
  - determining whether the sender is in a spam sender(s) store.
32. The method of claim 30, the challenge being at least one of an embedded code, a computational challenge, a human challenge and a micropayment request.
33. A method for responding to e-mail challenges, comprising:
  - receiving challenges to e-mail messages;
  - ordering the challenges based, at least in part, upon a message with fewer challenges processed before a message with more challenges;
  - processing the challenge of the message with fewer challenges; and,
  - sending a response to the challenge of the message with fewer challenges.
34. A data packet transmitted between two or more computer components that facilitates unsolicited e-mail detection, the data packet comprising:
  - a data field comprising information associated with a challenge, the challenge being based, at least in part, upon an associated probability that an e-mail message is spam.
35. A computer readable medium storing computer executable components of a system facilitating detection of unsolicited e-mail, comprising:
  - a mail classifier component that receives e-mail messages and determines an associated probability that the e-mail message is spam; and,
  - a challenge component that sends a challenge to a sender of an e-mail message having an associated probability greater than a first threshold.
36. A system facilitating detection of unsolicited e-mail, comprising:

means for determining an associated probability that an e-mail message is spam; and,

means for sending a challenge to a sender of an e-mail message having an associated probability greater than a first threshold.

5

10

15

20

25

30

35

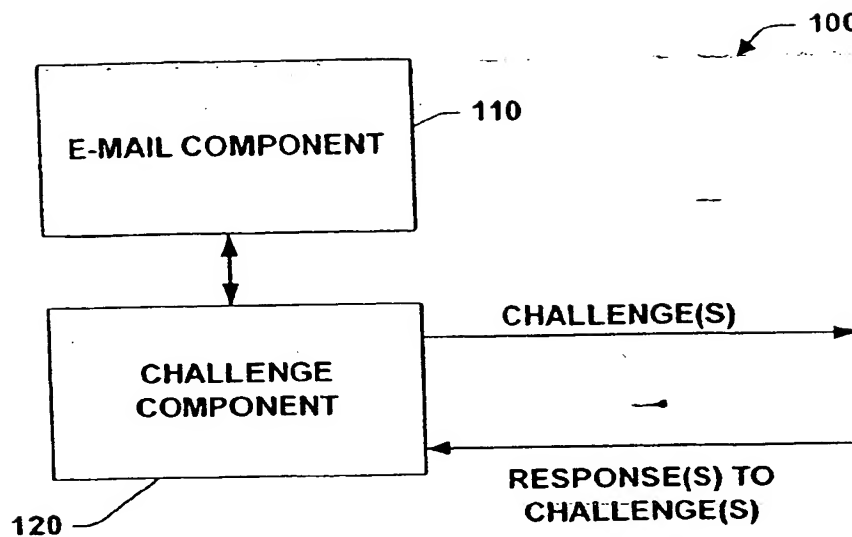
40

45

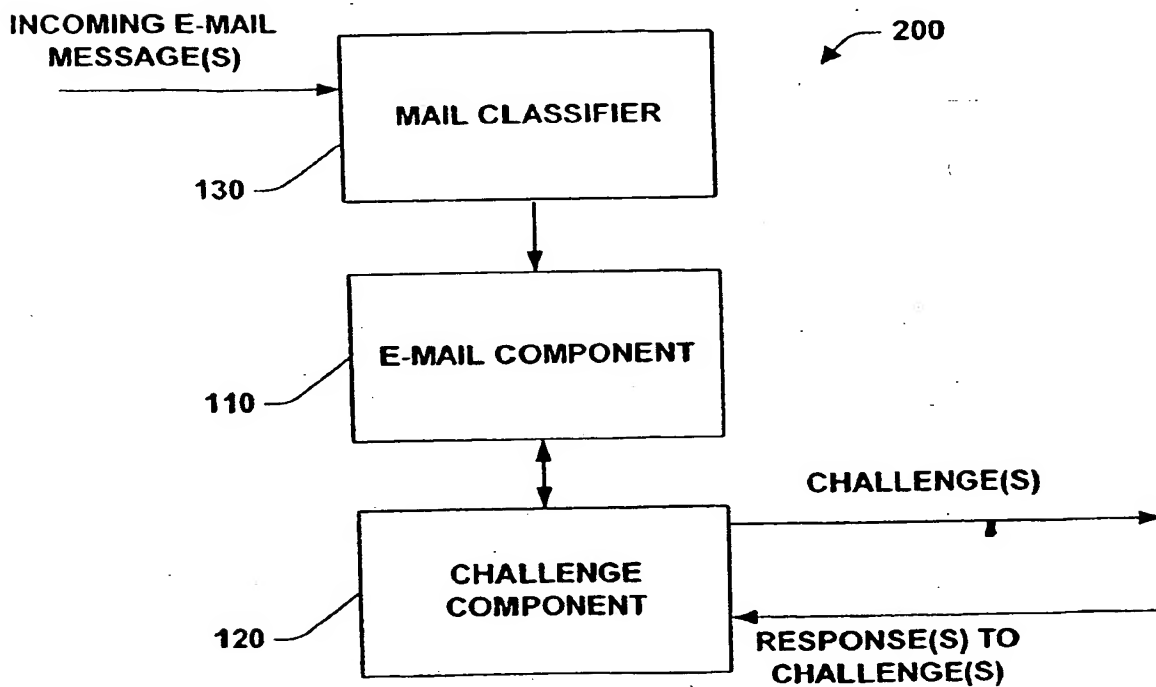
50

55

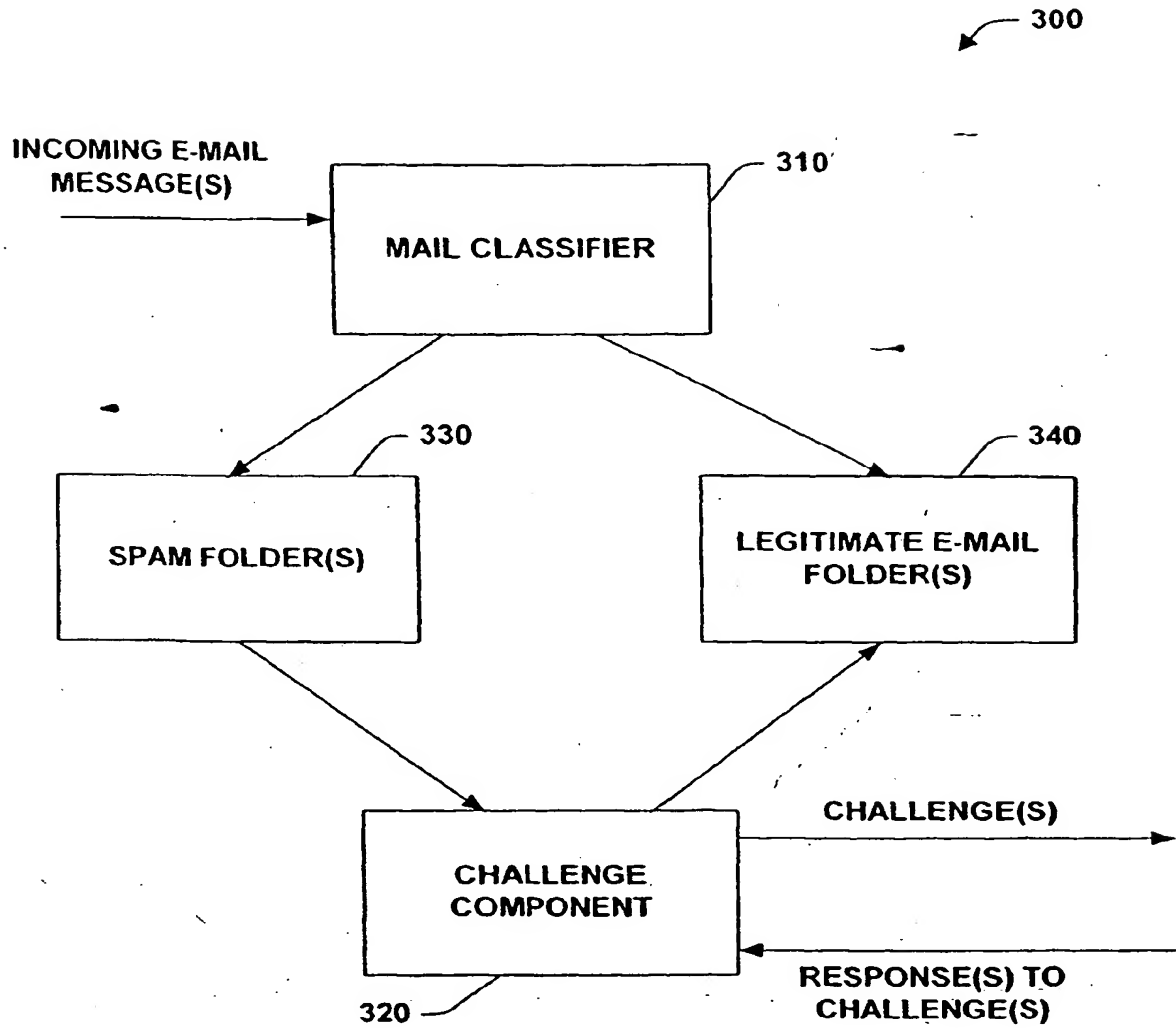
17



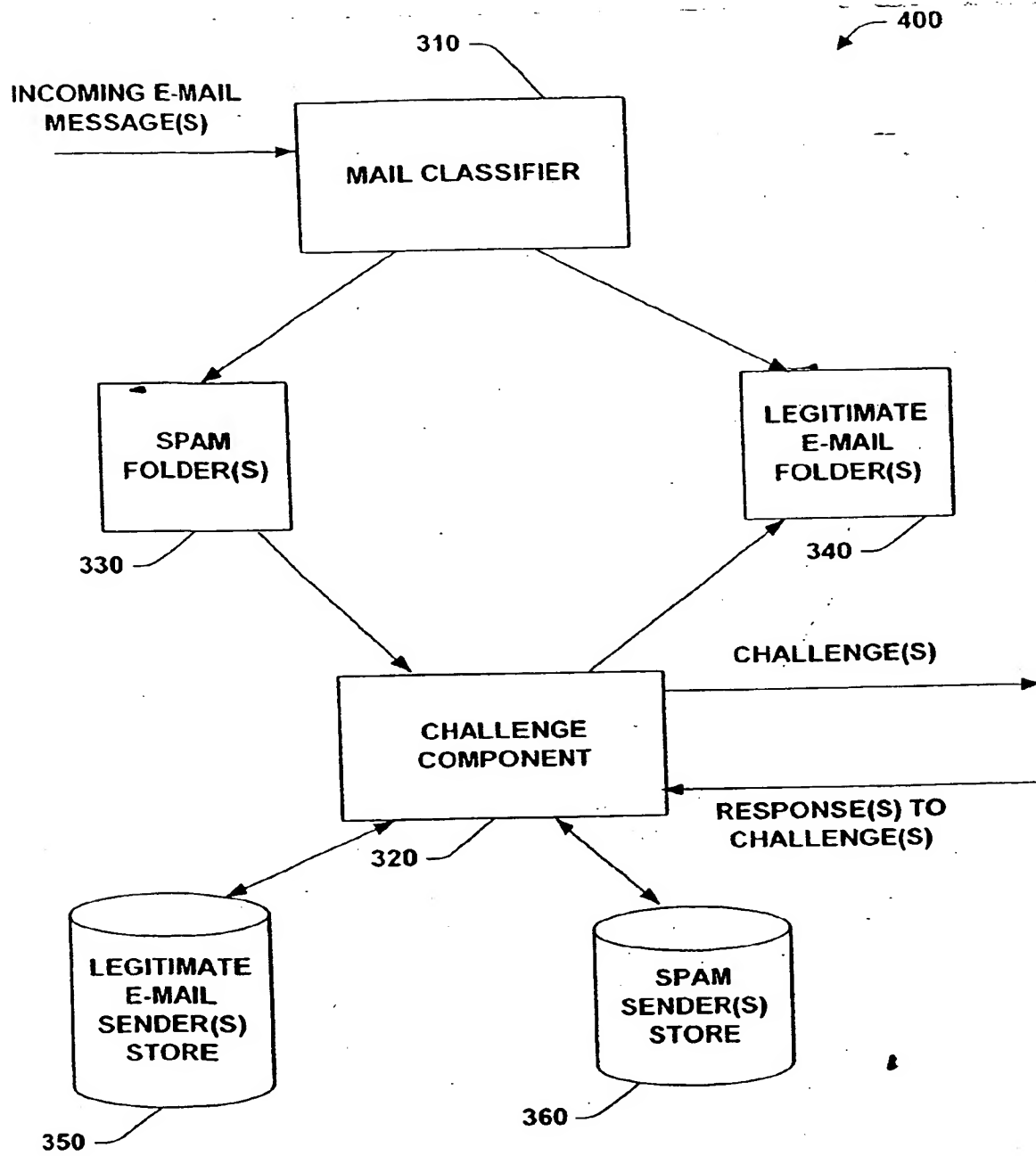
**FIG. 1**



**FIG. 2**



**FIG. 3**



**FIG. 4**

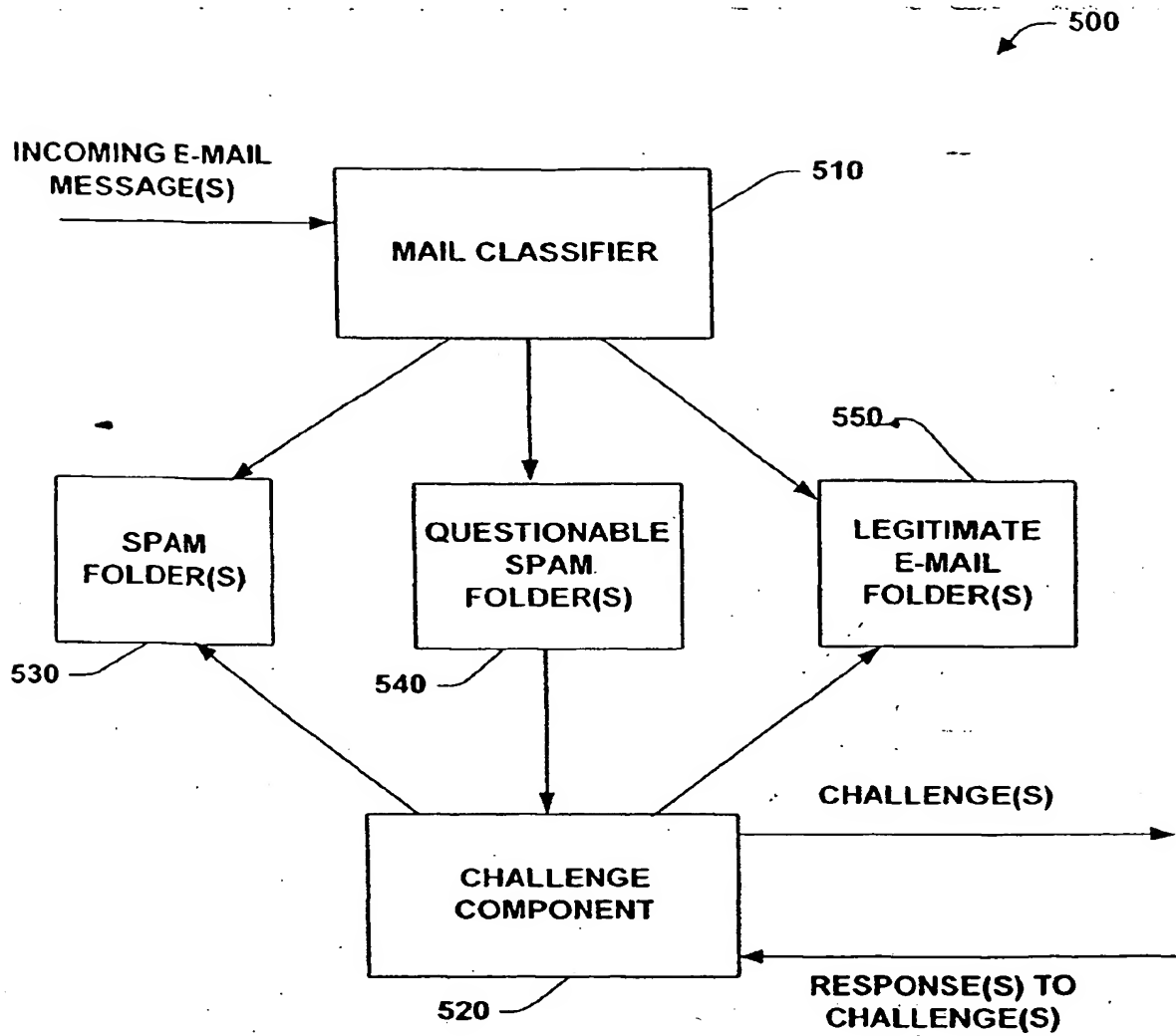
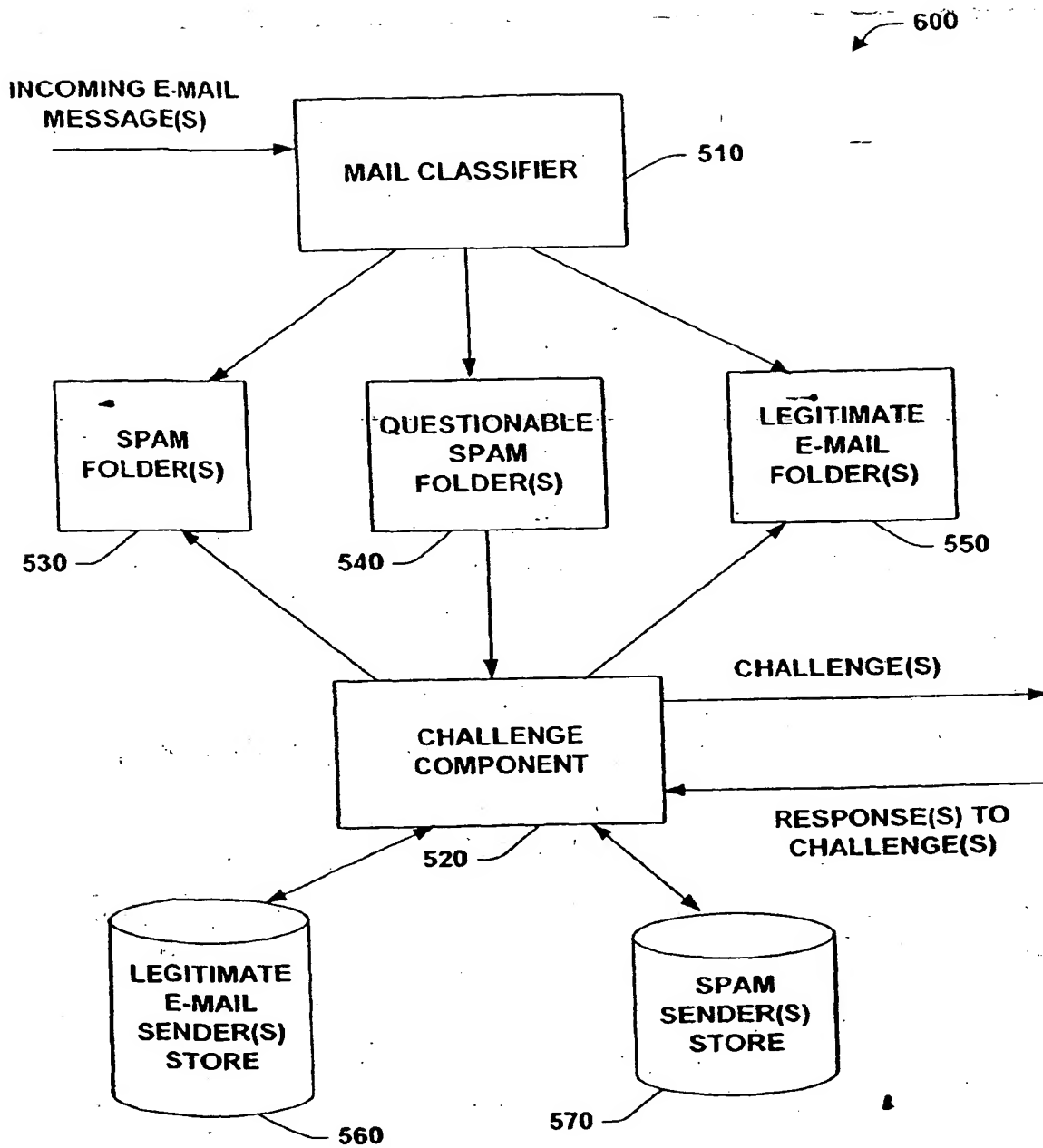
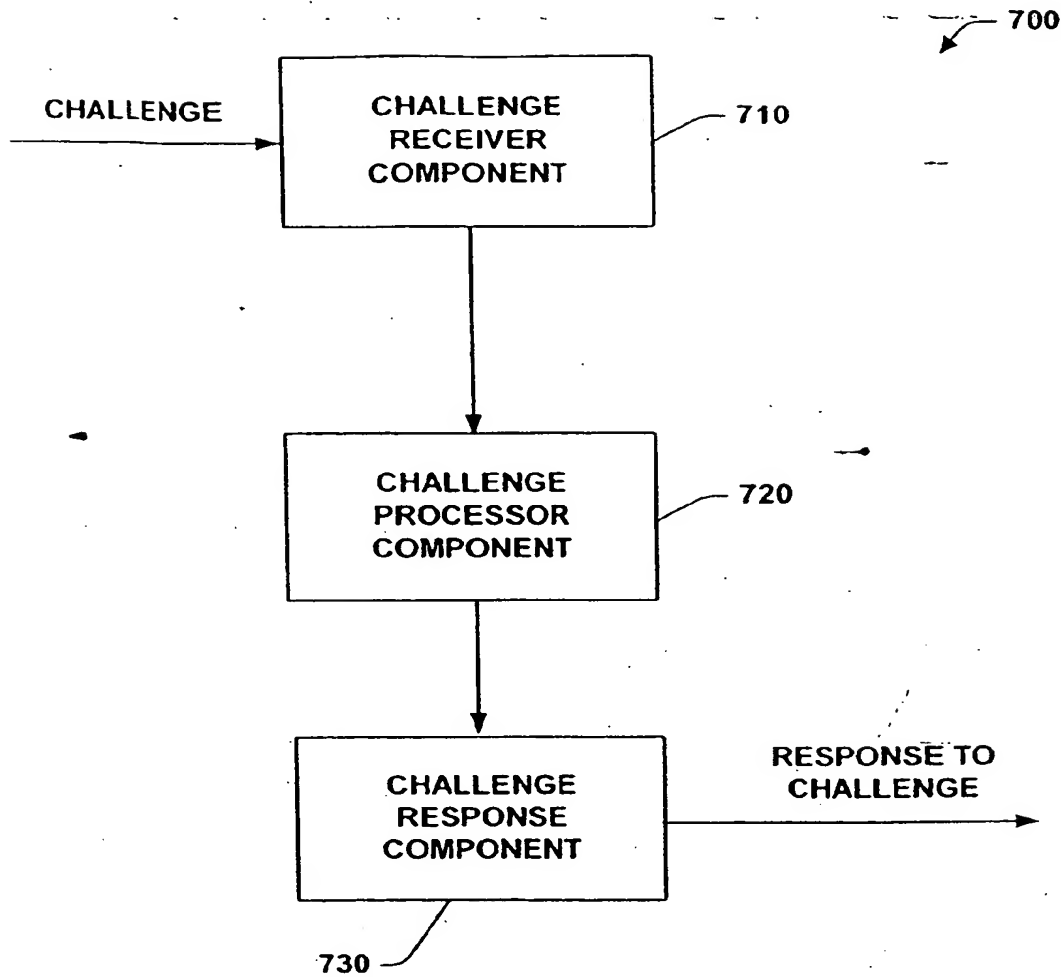


FIG. 5



**FIG. 6**



**FIG. 7**

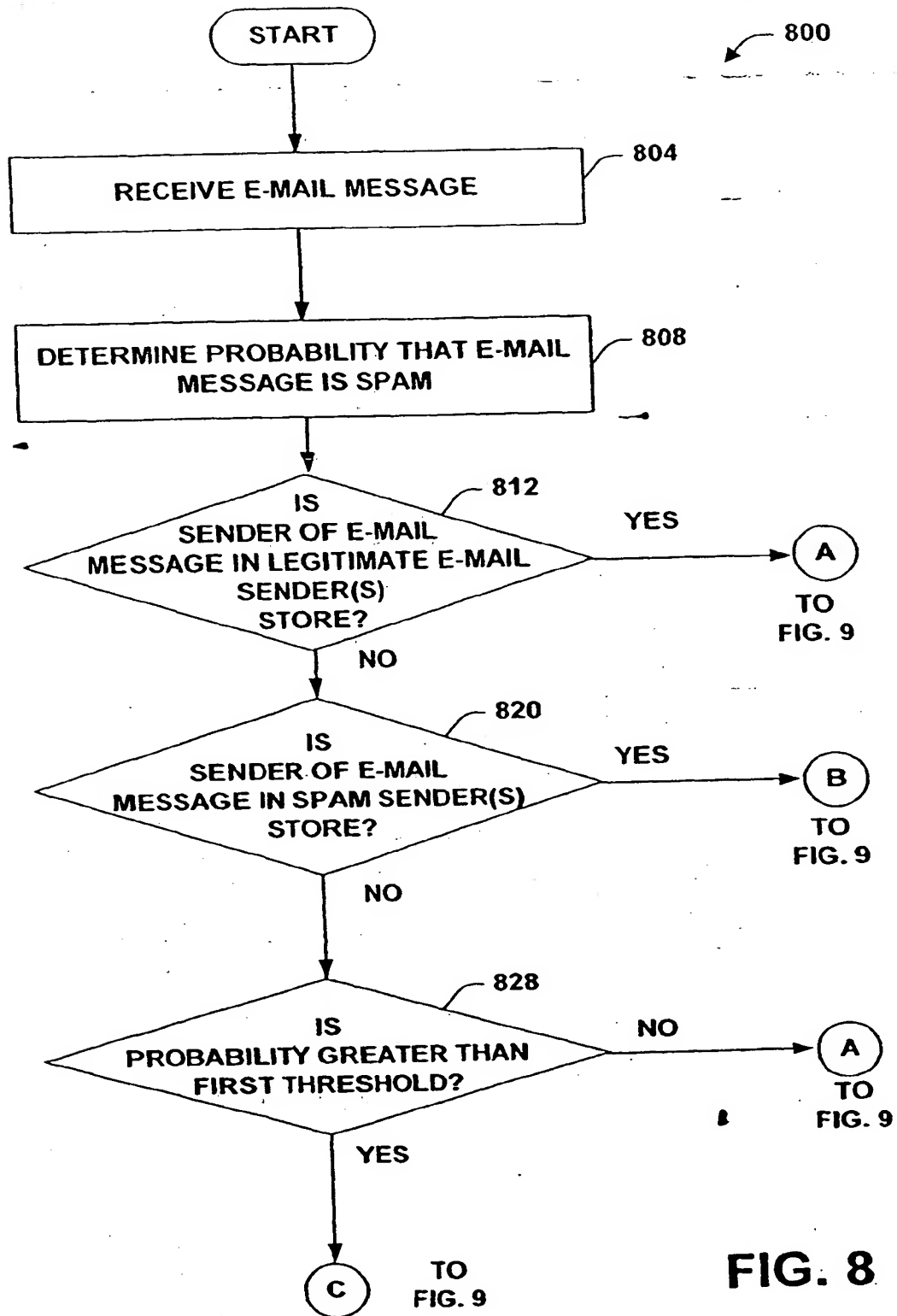
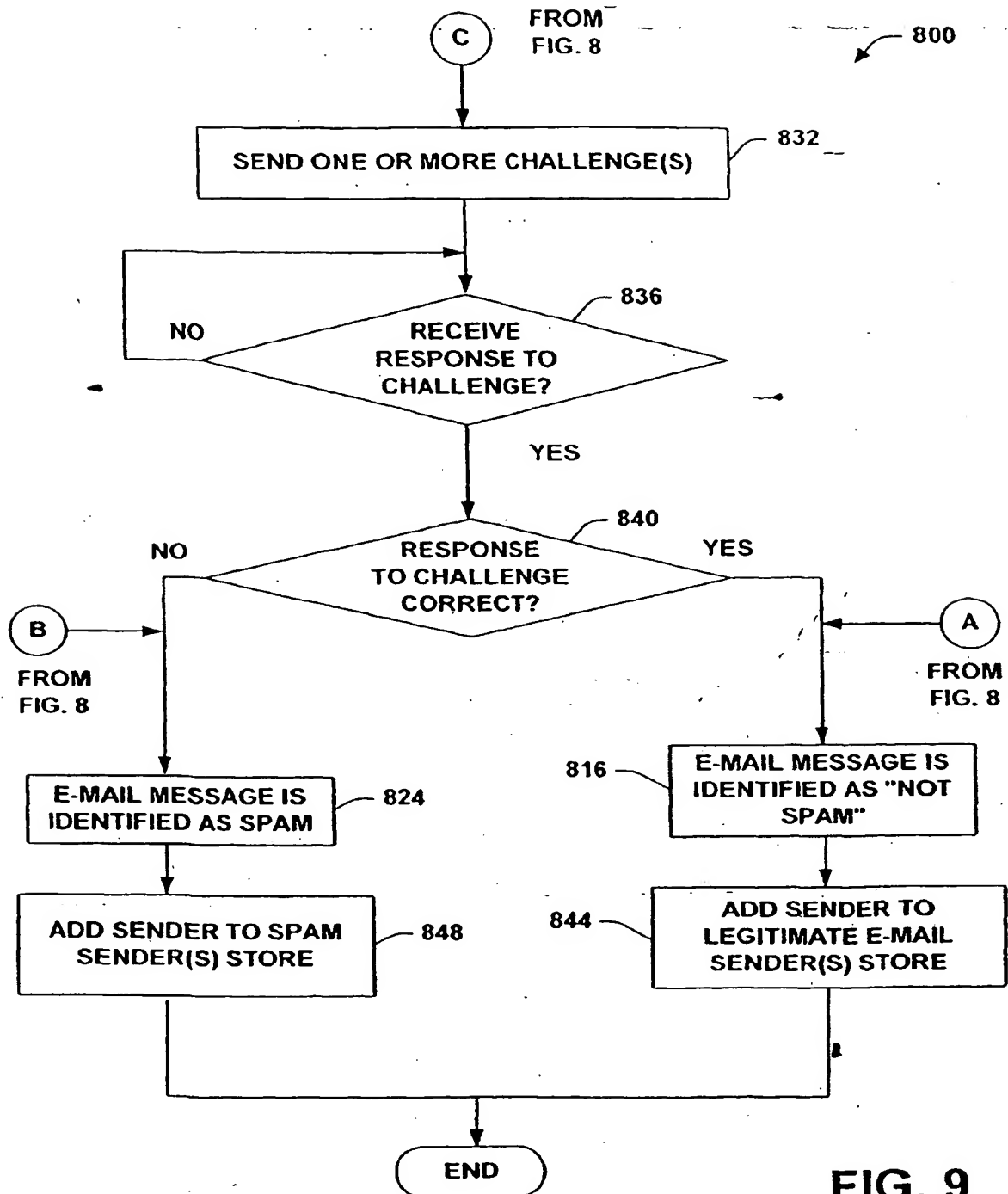
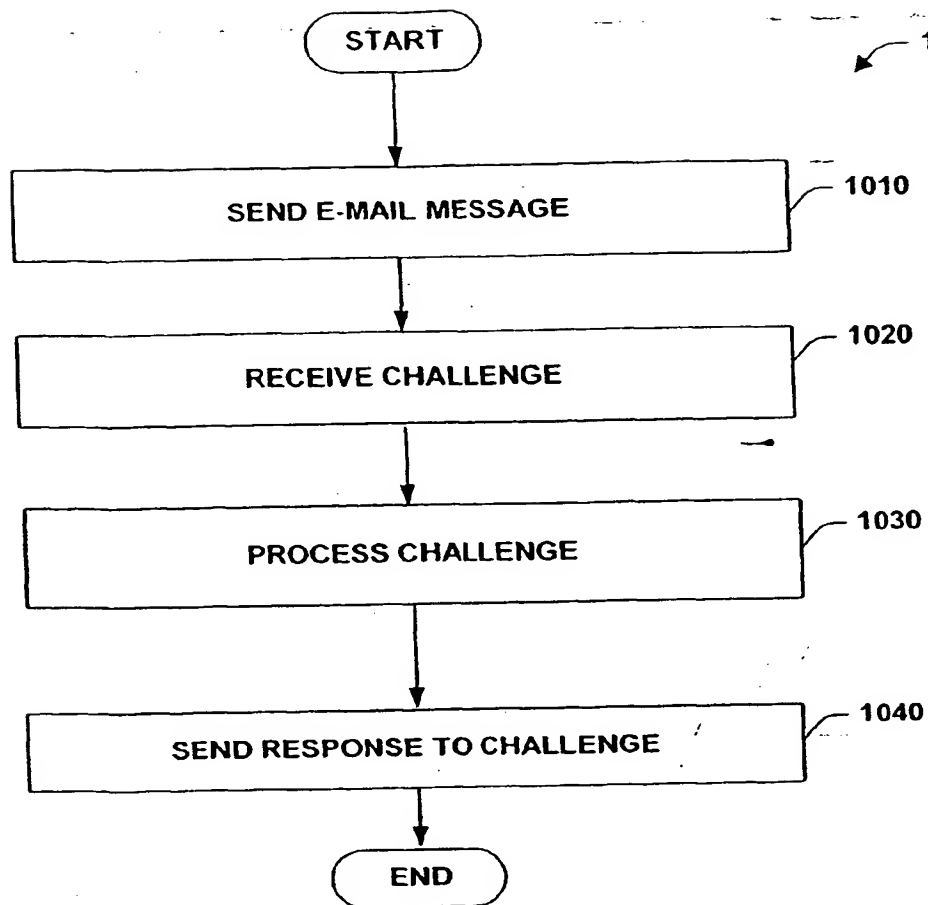
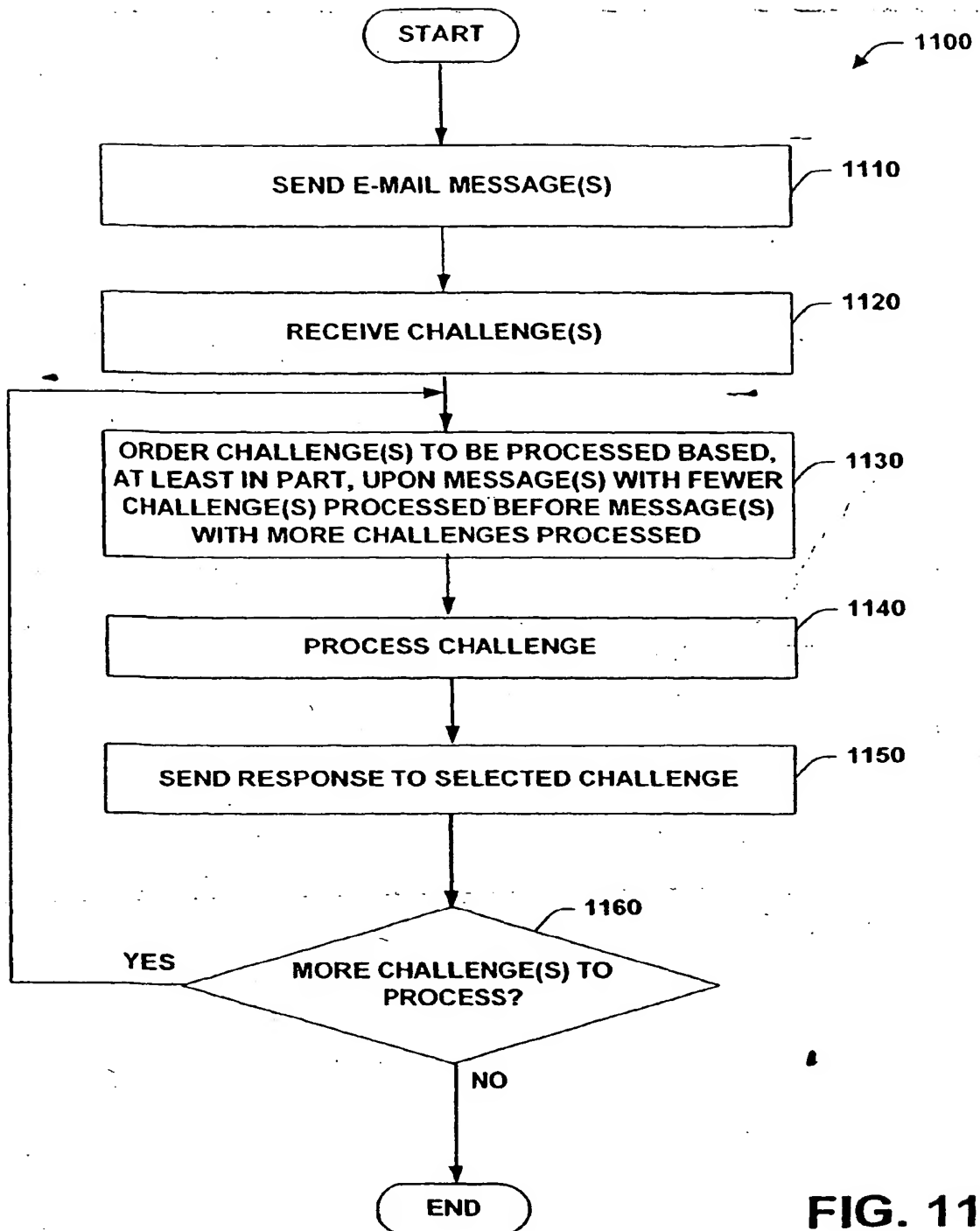


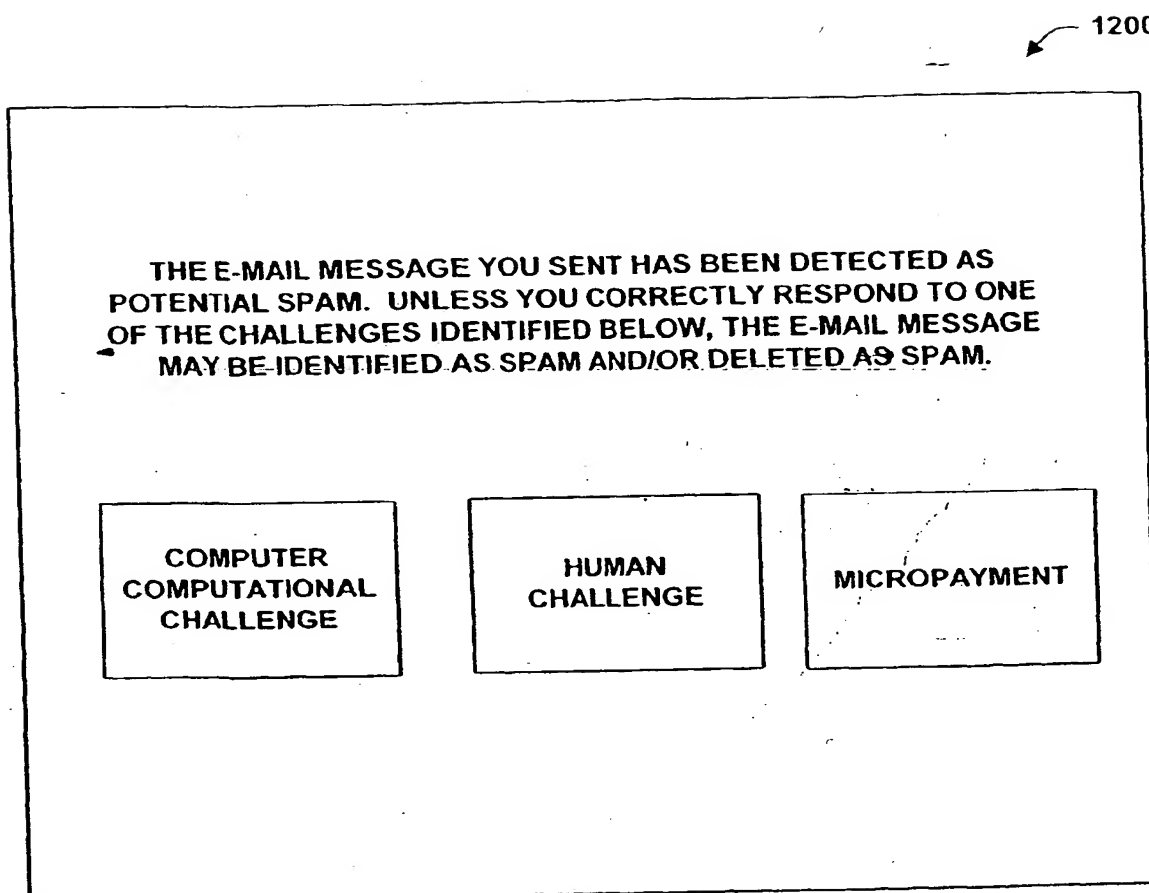
FIG. 8





**FIG. 10**

**FIG. 11**



**FIG. 12**

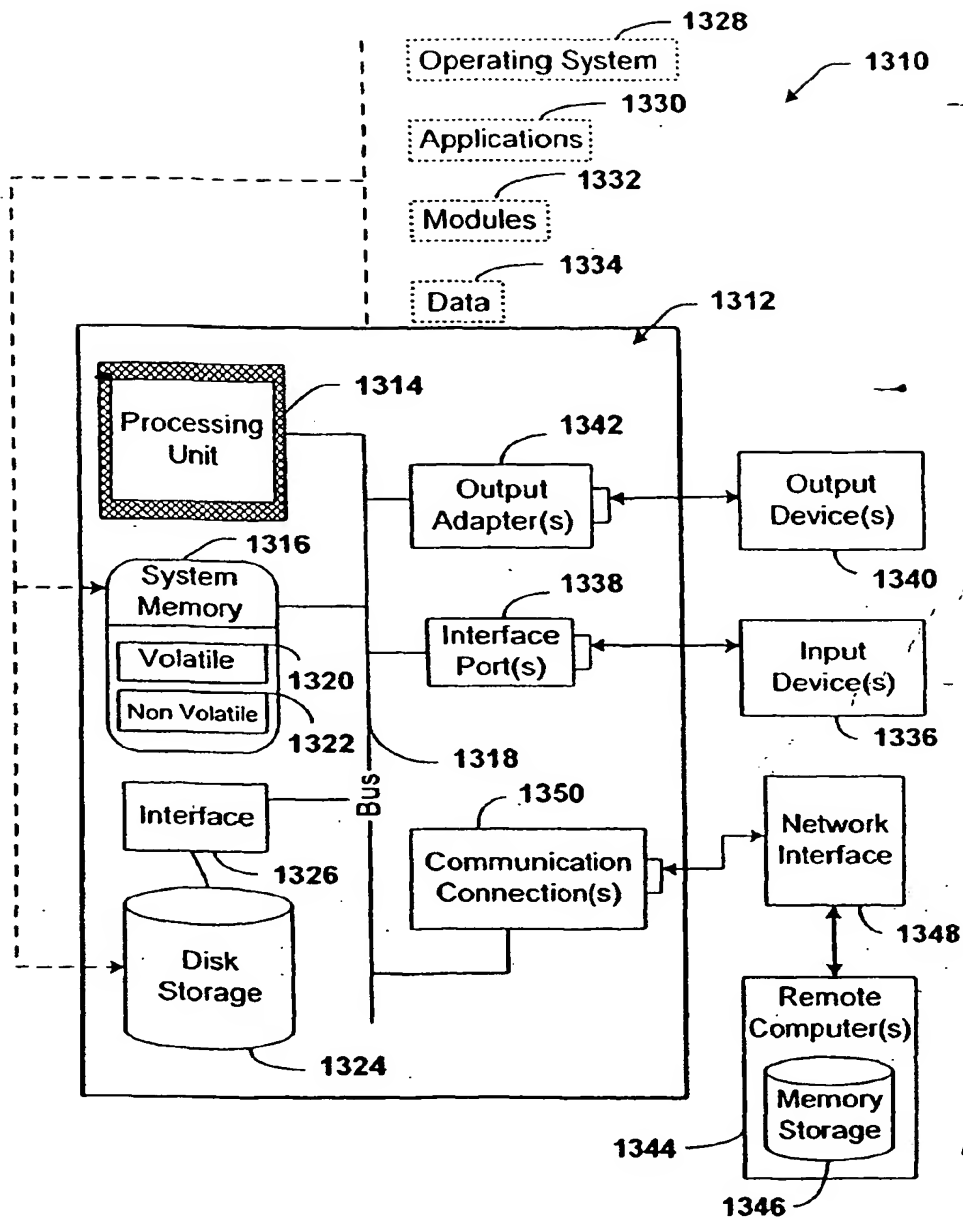
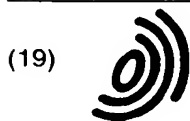


FIG. 13

**THIS PAGE BLANK (USPTO)**



Europäisches Patentamt  
European Patent Office  
Office européen des brevets



(11)

**EP 1 376 427 A3**

(12)

# EUROPEAN PATENT APPLICATION

(88) Date of publication A3:  
31.03.2004 Bulletin 2004/14

(51) Int Cl.7: **G06F 17/60**

(43) Date of publication A2:  
02.01.2004 Bulletin 2004/01

(21) Application number: **03006814.2**

(22) Date of filing: **26.03.2003**

(84) Designated Contracting States:  
**AT BE BG CH CY CZ DE DK EE ES FI FR GB GR  
HU IE IT LI LU MC NL PT RO SE SI SK TR**  
Designated Extension States:  
**AL LT LV MK**

(72) Inventors:  
• **Goodman, Joshua Theodore**  
**Redmond, Washington 98052 (US)**  
• **Rounthwaite, Robert L.**  
**Fall City, Washington 98024 (US)**

(30) Priority: **26.06.2002 US 180565**

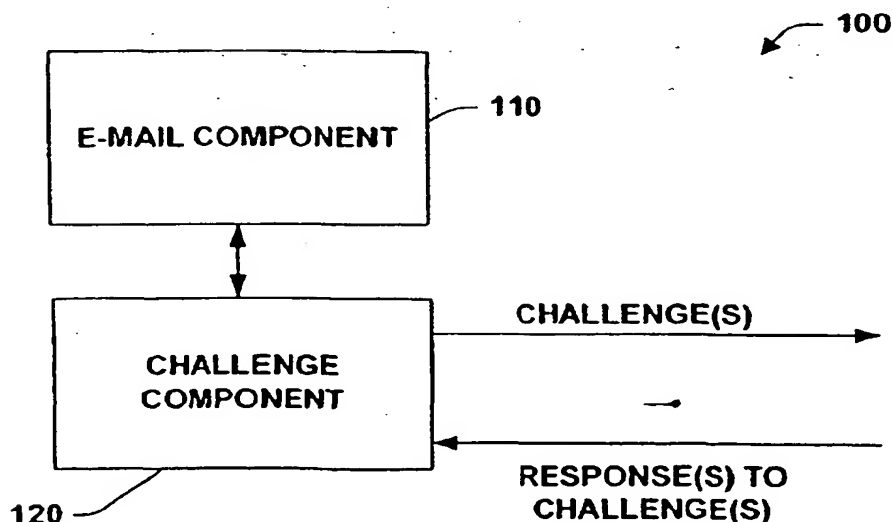
(74) Representative: **Grünecker, Kinkeldey,  
Stockmair & Schwanhäusser Anwaltssozietät**  
**Maximilianstrasse 58**  
**80538 München (DE)**

(71) Applicant: **MICROSOFT CORPORATION**  
**Redmond, Washington 98052-6399 (US)**

(54) **SPAM detector with challenges**

(57) A system and method facilitating detection of unsolicited e-mail message(s) with challenges is provided. The invention includes an e-mail component and a challenge component. The system can receive e-mail message(s) and associated probabilities that the e-mail message(s) are spam. Based, at least in part, upon the associated probability, the system can send a challenge

to a sender of an e-mail message. The challenge can be an embedded code, computational challenge, human challenge and/or micropayment request. Based, at least in part, upon a response to the challenge (or lack of response), the challenge component can modify the associated probability and/or delete the e-mail message.



**FIG. 1**

EP 1 376 427 A3



European Patent  
Office

## EUROPEAN SEARCH REPORT

Application Number

EP-03 00 6814

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.7)
D, Y	US 6 161 130 A (HECKERMAN DAVID E ET AL) 12 December 2000 (2000-12-12) * abstract * * column 1, line 10 - line 14 * * column 4, line 40 - line 49 * * column 8, line 40 - line 66 * * claim 1 *	1-36	G06F17/60
Y	WO 99 10817 A (COBB CHRISTOPHER ALAN) 4 March 1999 (1999-03-04) * abstract * * page 3, line 4 - line 17 * * page 7, line 19 - page 8, line 10 * * page 16, line 19 - line 23 * * page 11, line 18 - line 21 * * page 28, line 27 - line 30 * * figure 7 *	1-36	
Y	US 6 112 227 A (HEINER JEFFREY NELSON) 29 August 2000 (2000-08-29) * abstract * * column 2, line 1 - line 10 * * column 3, line 58 - column 4, line 1 * * column 4, line 15 - line 30 * * figure 2 *	1-36	TECHNICAL FIELDS SEARCHED (Int.Cl.7)
A	JULIAN BYRNE: "MY Spamblock" NEWSGROUP CITATION, 'Online! 19 January 1997 (1997-01-19), XP002267503 news.admin.net-abuse.email Retrieved from the Internet: <URL:http://www.google.com/groups?hl=en&lr= =&ie=UTF-8&oe=UTF-8&selm=32E1A4FD.41C6%40a ny.where&rnum=1> 'retrieved on 2004-01-20! * the whole document *	1-36	G06F
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 21 January 2004	Examiner Rossier, T
CATEGORY OF CITED DOCUMENTS X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document		T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons S : member of the same patent family, corresponding document	

EPO FORM 1503 03 82 (Pst/C01)



European Patent  
Office

# EUROPEAN SEARCH REPORT

Application Number

EP 03 00 6814

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.7)
A	<p>DAVID SKOLL: "How to make sure a human is sending you mail"</p> <p>NEWSGROUP CITATION, 'Online!'</p> <p>17 November 1997 (1997-11-17), XP002267504</p> <p>news.admin.net-abuse.usenet</p> <p>Retrieved from the Internet:</p> <p>&lt;URL:http://groups.google.ca/groups?hl=en&amp;lr=&amp;ie=UTF-8&amp;oe=UTF-8&amp;selm=561uge%246on%40bertrand.ccs.carleton.ca&gt;</p> <p>'retrieved on 2004-01-20!'</p> <p>* the whole document *</p> <p>-----</p>	1-36	
			TECHNICAL FIELDS SEARCHED (Int.Cl.7)
The present search report has been drawn up for all claims			
Place of search		Date of completion of the search	Examiner
THE HAGUE		21 January 2004	Rossier, T
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone</p> <p>Y : particularly relevant if combined with another document of the same category</p> <p>A : technological background</p> <p>O : non-written disclosure</p> <p>P : intermediate document</p> <p>T : theory or principle underlying the invention</p> <p>E : earlier patent document, but published on, or after the filing date</p> <p>D : document cited in the application</p> <p>L : document cited for other reasons</p> <p>a : member of the same patent family, corresponding document</p>			

EPO FORM 1501 03 92 (P/AC/01)

# ANNEX TO THE EUROPEAN SEARCH REPORT ON EUROPEAN PATENT APPLICATION NO.

EP 03 00 6814

This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report. The members are as contained in the European Patent Office EDP file on  
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

21-01-2004

Patent document cited in search report		Publication date		Patent family member(s)	Publication date
US 6161130	A	12-12-2000	EP	1090368 A1	11-04-2001
			WO	9967731 A1	29-12-1999
WO 9910817	A	04-03-1999	US	6199102 B1	06-03-2001
			WO	9910817 A1	04-03-1999
US 6112227	A	29-08-2000	NONE		

EPO FORM 1/0459

For more details about this annex : see Official Journal of the European Patent Office, No. 12/82